wnloaded at #61431990 on January 7, 202

On the impossibility of predicting the behavior of rational agents

Dean P. Foster[†] and H. Peyton Young^{‡§}

[†]Department of Statistics, Wharton School, University of Pennsylvania, Philadelphia, PA 19104; and [‡]Department of Economics, The Johns Hopkins University, Baltimore, MD 21218-2685

Edited by Richard D. McKelvey, California Institute of Technology, Pasadena, CA, and approved August 1, 2001 (received for review November 8, 2000)

A foundational assumption in economics is that people are rational: they choose optimal plans of action given their predictions about future states of the world. In games of strategy this means that each player's strategy should be optimal given his or her prediction of the opponents' strategies. We demonstrate that there is an inherent tension between rationality and prediction when players are uncertain about their opponents' payoff functions. Specifically, there are games in which it is impossible for perfectly rational players to learn to predict the future behavior of their opponents (even approximately) no matter what learning rule they use. The reason is that in trying to predict the next-period behavior of an opponent, a rational player must take an action this period that the opponent can observe. This observation may cause the opponent to alter his next-period behavior, thus invalidating the first player's prediction. The resulting feedback loop has the property that, a positive fraction of the time, the predicted probability of some action next period differs substantially from the actual probability with which the action is going to occur. We conclude that there are strategic situations in which it is impossible in principle for perfectly rational agents to learn to predict the future behavior of other perfectly rational agents based solely on their observed actions.

Rationality vs. Predictability

E conomists often assume that people are *rational*: they maximize their expected payoffs given their beliefs about future
states of the world. This hypothesis plays a crucial role in game
theory, where each player is assumed to choose an optimal
strategy given his belief about the strategies of his opponents. In
this setting, a belief amounts to a forecast or prediction of the
opponents' future behavior, that is, of the probability with which
the opponents will take various actions. The prediction is *good*if the forecasted probabilities are close to the actual probabilities. Together prediction and rationality justify the central
solution concept of the theory. Namely, if each player correctly
predicts the opponents' strategies and if each chooses an optimal
strategy given his prediction, then the strategies form a Nash
equilibrium of the repeated game. But under what circumstances
will rational players actually learn to predict the behavior of
others starting from out-of-equilibrium conditions?

In this article we show that there are very simple games of incomplete information such that players almost never learn to predict their opponents' behavior even approximately, and they almost never come close to playing a Nash equilibrium. This impossibility result and its proof builds on the existing literature on learning in repeated games (1–8); for other critiques of Bayesian learning in economic environments see refs. 9–11 and 19. The present contribution demonstrates the incompatibility between rationality and prediction without placing any restrictions on the players' prior beliefs, their learning rules, or the degree to which they are forward-looking.

An Example. We begin by illustrating the problem in a concrete case. Consider two individuals, A and B, who are playing the game of matching pennies. Simultaneously each turns a penny

face up or face down. If the pennies match (both are heads or both are tails), then B buys a prize for A; if they do not match, A buys a prize for B. Assume first that the prize is one dollar and that the utility of both players is linear in money. Then the game has a unique Nash equilibrium in which each player randomizes by choosing heads (H) and tails (T) with equal probability. If both adopt this strategy, then each is optimizing given the strategy of the other. Moreover, although neither can predict the *realized action* of the opponent in any given period, each can predict his strategy, namely, the *probabilities* with which the actions will be taken. In this case no tension exists between rationality and prediction because the game has a unique equilibrium, and the players know what it is.

Now change the situation by assuming that if both players choose H, then B buys an ice cream cone for A, whereas if both choose T then B buys A a milk shake. Similarly, if A chooses H and B chooses T then A buys B a coke, whereas if the opposite occurs then A buys B a bag of chips. Assume that the game is played once each day, the players' tastes do not change from one day to the next, and they have a fixed positive utility for each of the prizes and also for money. Unlike the previous situation, this is a game of incomplete information in which neither player knows the other's payoffs.

		B's	action	B's action			
		Н	T	Н	T		
A's action	Н	eat cone	buy coke	buy cone	drink coke		
	T	buy chips	drink shake	eat chips	buy shake		

Outcomes for A Outcomes for B

For expositional simplicity assume first that the players are myopic, that is, they do not worry about the effect of their actions on the future course of the game. Imagine that the following sequence of actions has occurred over the first ten periods.

Period	1	2	3	4	5	6	7	8	9	10	11
A:	Η	T	T	Η	Η	Η	T	Η	T	Η	?
B:	Τ	Η	T	Η	T	Η	T	Η	T	Η	?

The immediate problem for each player is to predict the *intention* of the opponent in period 11 and to choose an optimal response. The opponent's intention might be to play H for sure, T for sure, or to randomize with probability p for H and 1 - p for T. If the

This paper was submitted directly (Track II) to the PNAS office.

Abbreviations: H, heads; T, tails; i.i.d., Independent and identically distributed.

[§]To whom reprint requests should be addressed. E-mail: pyoung@jhu.edu.

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

opponent's intention is to randomize, then obviously one cannot predict his realized action, but it does not seem too much to ask that one predict the approximate probability with which he intends to play each action. We claim, however, that this is essentially impossible.

To see why, let's put ourselves in A's shoes. The behavior of B suggests an alternating pattern, perhaps leading us to predict that B will play T next period. Because we are rational, we will (given our prediction) play T for sure next period. But if B is a good predictor, then she must be able to predict that with high probability we are in fact going to play T next period. This prediction by B leads her to play H next period, thus falsifying our original prediction that she is about to play T.

The point is that if either side makes a prediction that leads them to play H or T for sure, the other side must predict that they are going to do so with high probability, which means that they too will choose H or T for sure. But there is no pair of predictions such that both are approximately correct and the optimal responses are H or T for sure. It follows that, for both players to be good predictors of the opponent's next-period behavior, at least one of them must be intending to play a mixed strategy next period, and the other must predict this.

Suppose, for example, that player B intends to play a mixed strategy in period 11. Because B is rational, she only plays a mixed strategy if she is *exactly indifferent* between playing H and T given her predictions about A. (If there is a slight difference in payoff between the two actions, strict rationality requires that the one with higher payoff be chosen exclusively.) Now B's predictions about A's behavior in the 11th period are based on the observed history of play in the first 10 periods. Let's say that the particular history given above leads B to predict that A will play H with a probability of 0.127. Because B intends to play mixed, it must be the case that B's expected utility from playing H or T is identical given B's utility function $u_{\rm B}$ for the various outcomes. In other words, it must be that

.127 $u_{\rm B}$ (buy cone) + .873 $u_{\rm B}$ (eat chips) = .127 $u_{\rm B}$ (drink coke) + .873 $u_{\rm B}$ (buy shake).

But there is no reason to think that B's utilities actually do satisfy this equation exactly. More precisely, let us suppose that B's utility for each outcome could be any real number within a certain interval, and that B's actual utility (B's type) is the result of a random draw from among these possible values. (The draw occurs once and for all before the game begins.) Following Jordan (2), we claim that the probability is zero that the above equation will be satisfied. The reason is that there is only a finite number of distinct predictions that B could make at this point in time, because B's prediction can only be based on A's (and B's) previous observed behavior together with B's initial beliefs. Because this argument holds for every period, the probability is zero that B will ever be indifferent. From this and the preceding argument it follows that in any given period, one or both players must be making a bad prediction. Moreover, they cannot be playing a Nash equilibrium in any given period (or even close to a Nash equilibrium), because this would require them to play mixed strategies, which means that both must be indifferent.

Jordan (2) was the first to employ this kind of argument to show that myopic players effectively cannot learn mixed equilibria no matter what their beliefs are. Moreover, as we have just seen, the same argument shows that at least one of them cannot learn to predict the behavior of the other. The limitation of Jordan's result is that it assumes players are completely myopic. Forward-looking behavior allows for a richer repertoire of learning strategies and more time to detect complex patterns in the behavior of one's opponent. Nevertheless, the incompatibil-

ity between rationality and prediction continues to hold even in this case, as we shall show below.

A second closely related body of work is provided by Nachbar (6–8). He was the first to argue that there is a fundamental tension between prediction and rationality in the context of Bayesian learning even when players are forward-looking. Nachbar's critique was prompted by an earlier paper by Kalai and Lehrer (4), which laid out conditions under which Bayesian rational players would in fact be able to learn to predict the behavior of their opponents. Suppose that each player begins the game with a prior belief over the possible repeated game strategies that his opponents might use. Kalai and Lehrer show that, if these prior beliefs contain a "grain of truth," that is, they put positive probability (however small) on the *actual* repeated game strategies of the opponents, then players learn to predict with probability one.

As Nachbar points out, however, the grain-of-truth condition may be very difficult to satisfy in practice. To illustrate, consider the preceding example and suppose that the players are perfectly myopic. Then the unique equilibrium of the repeated game is for A to play H with some fixed probability p^* each period and for B to play H with some fixed probability q^* each period. These values are not known to the players because p^* depends on B's payoffs, whereas q^* depends on A's payoffs. Can they be learned through Bayesian updating of a diffuse prior? Suppose that each player begins with a belief that the other is playing an i.i.d. strategy with an unknown parameter (the probability of playing H), where the beliefs have full support on the interval [0, 1]. In any given period, the players almost surely will have updated beliefs that lead them to play H or T with a probability of 1 in that period, because the expected payoffs from H and T are not exactly equal. However, their updated beliefs lead them to predict that their opponent is almost surely going to play a mixed strategy next period. Thus their predictions almost certainly are not close to their actual strategies. Furthermore, as the game proceeds, rationality causes them to play H for sure in some periods and T for sure in others. Hence their actual strategies are not i.i.d. and not in the support of their beliefs. More generally, Nachbar (6-8) argues that in games such as this it is difficult to identify any plausible family of beliefs such that the players' best response strategies are in the support of their beliefs [another paper in the same general spirit is provided by R. I. Miller and C. W. Sanchirico (12)].

In this paper we are agnostic about whether or not the players are Bayesian and what the structure of their priors might be. Instead we show that no matter how players use the information revealed by repeated play, they will fail to learn to predict the opponents' behavior in some kinds of games.

Before turning to a precise statement of our result, we should point out that it is *prediction by the players* that is problematical; to an observer the average behavior of the players may exhibit empirical regularities. For example, it could be that the cumulative frequency distribution of play approaches a Nash equilibrium of the game. In fact this will be the case for fictitious play in which each player uses the empirical distribution of the opponent's play up through a given period to predict his nextperiod behavior and then chooses a best response given that prediction. In games such as matching pennies, this simple learning rule induces long run average behavior that converges to the mixed Nash equilibrium of the game (13, 14). There are other models in which a player's average behavior mimics Nash equilibrium from the observer's standpoint (15, 16); in fact Nash himself proposed such an interpretation (17). But this does not imply that individual players ever play Nash equilibrium strategies or that they learn to predict.

The Learning Model. We now describe our impossibility result in detail. Consider an n-person game G with finite action space X = G

loaded at #61431990 on January 7, 2020

 Π X_i and utility functions u_i : $X \to R$. We shall assume that the payoffs take the form $u_i(x) = u_i^0(x) + \omega_i(x)$, where the $u_i^0(x)$ are payoffs in a benchmark game G^0 , and the $\omega_i(x)$ are i.i.d. random variables drawn from a continuous density $\nu(\omega)$, the support of which is the interval $I_{\lambda} = [-\lambda/2, \lambda/2]$. The parameter $\lambda > 0$ is the range of uncertainty in the payoffs. We shall assume that the distribution of payoffs is common knowledge, but the realized payoff $u_i(x)$ is known only to player i. Errors are drawn once only before play begins, and the resulting one-shot game (called the ν -perturbation of G^0) is played infinitely often.

Each player takes an action once in each time period $t = 1, 2, 3, \ldots$ The *outcome* in period t is an *n*-tuple of actions $x^t \in X$, where x_i^t is the action taken by *i* in period *t*. A *state* of the process at time *t* is a history of play up to *t*, that is, a sequence of outcomes $h^t = (x^1, x^2, \ldots x^t)$. Let h^0 represent the null history, H^t the set of all length-*t* histories, and $H = \bigcup_t H^t$ the set of all finite histories, i.e., the set of all states. A realization of the process will be denoted by h, and the set of realizations (i.e., the set of infinite histories) will be denoted by H^∞ . Histories are observed publicly, that is, there is perfect monitoring.

The discounted payoff to player i from a realization $h = (x^1, x^2, \dots x^t, \dots)$ is

$$U_i(h) = (1 - \delta_i) \sum_{t=1}^{\infty} \delta_i^{t-1} u_i(x^t),$$

where δ_i is i's discount factor, $0 \le \delta_i < 1$ (if $\delta_i = 0$, $U_i(h) = u_i(x^1)$). Let Δ_i denote the set of probability distributions over X_i . Let $\Delta = \Pi_i \Delta_j$ denote the product set of mixtures, and let $\Delta_{-i} = \Pi_{j \ne i} \Delta_j$ be the product set of mixtures by i's opponents. A behavioral strategy for player i specifies a conditional probability distribution over i's actions in each period conditional on the state in the previous period. Thus we can represent i's strategy by a function $q_i^t = g_i(h^{t-1}) \in \Delta_i$, where $q_i^t(x_i)$ is the probability that i plays x_i in period t given that t^{t-1} is the state in period t-1. This is of course a function of i's realized utility function u_i , but we shall not write this dependence explicitly.

A prior belief of player i is a probability distribution over all possible combinations of the opponent's strategies. We can decompose any such belief into one-step-ahead forecasts of the opponent's behavior conditional on each possible state. Thus, if h^{t-1} is the state at time t-1, i's forecast about the behavior of her opponents in period t can be represented by a probability distribution $p_{-i}^t = f_i(h^{t-1}) \in \Delta_{-i}$, where $p_{-i}^t(x_{-i})$ is the probability that i assigns to the others playing the combination x_{-i} in period t. The function $f_i : H \to \Delta_{-i}$ will be called i's forecasting function. Given any vector of forecasting functions $f = (f_1, f_2, \ldots f_n)$, one for each player, there exists a set of prior beliefs such that the f_i describe the one-step-ahead forecasts of players with these beliefs [see Kalai and Lehrer (4)].

Consider the situation just after the players have been informed privately of their realized payoff functions u_i . Because of the independence of the draws among players, no one knows anything he did not already know about the others' payoffs, and this fact is common knowledge. This has an implication for the forecasting functions. Namely, at the beginning of each period t, i knows that j's information consists solely of the publicly observed history h^{t-1} and j's own payoff function u_j . Player j's behavior cannot be conditioned on information that *j* does not have (namely u_{-i}), and player i's forecast of j's behavior cannot be conditioned on information that *i* does not have (namely, u_{-i}). Thus i's forecast $[f_i(h^{t-1})]_j$ about j's behavior in each period t does not depend on the realization of the values u_k for every k, including k = i, j. It follows that the functions f_i do not depend on the realized payoff functions $u_i(\cdot)$, although they may depend on ν . Another way of saying this is that the beliefs must be consistent with the players' a priori knowledge of the information structure.

Following Jordan (2), we shall say that a *learning process* is a pair $(f, g) = (f_1, \ldots, f_n, g_1, \ldots, g_n)$, where $f_i : H \to \Delta_{-i}$ and $g_i : H \to \Delta_i$ for each player i. Given a realization of the process h, we shall denote player i's forecast in period t by $p_{-i}^t(h) = f_i(h^{t-1})$, and i's behavioral strategy in period t by $q_i^t(h) = g_i(h^{t-1})$.

The pair (f_i, g_i) induces a probability measure on the set of all realizations H^{∞} . Similarly, for every state h^{t-1} , f_i and g_i induce a conditional probability distribution on all continuations of h^{t-1} . Denote this conditional distribution by $\mu_i(f_i, g_i | h^{t-1})$. We say that individual i is *rational* if, for every h^{t-1} , i's conditional strategy $g_i(\cdot|h^{t-1})$ optimizes i's expected utility from time t on, given i's conditional forecast $f_i(\cdot|h^{t-1})$. (This is also known as *sequential rationality*.) Specifically, for every alternative choice of strategy $g_i'(\cdot|h^{t-1})$,

$$\int U_i(h)d\mu_i(f_i,g_i|h^{t-1}) \geq \int U_i(h)d\mu_i(f_i,g_i'|h^{t-1}).$$

Prediction. Intuitively, player i learns to predict the behavior of his opponent(s) if i's forecast of their next-period behavior comes closer and closer to their actual next-period strategies. This idea may be formalized as follows. Consider a learning process (f, g), and let $\mu(g)$ denote the probability measure induced on H^{∞} by the strategies $g = (g_1, g_2, \ldots, g_n)$. We say that player i learns to predict if the mean square error of i's next-period predictions goes to zero over almost all histories of play. In other words, for $\mu(g)$ -almost all realizations h

$$\lim_{T \to \infty} \sum_{t=1}^{T} |p_{-i}^{t}(h) - q_{-i}^{t}(h)|^{2}/T = 0.$$
 [1]

Similarly, we shall say that player *i never learns to predict* if the subset of histories for which Eq. 1 holds has μ -measure zero. Note that this condition permits players to make bad forecasts from time to time, provided they do not occur too often.

An Impossibility Theorem. We now demonstrate a class of repeated games such that, with probability one, some player never learns to predict his opponent's behavior, and this holds for all prior beliefs. Because our result holds for all beliefs, it must hold for beliefs that are in some sense best possible. A reasonable candidate for "best possible beliefs" are rational expectations beliefs. These have the property that, at every point in time, each player's prediction of his opponent's future behavior is conditioned correctly on the posterior distribution of payoff types revealed by play so far. Jordan (1, 3) shows that these posterior distributions converge to the set of Nash equilibria of the game (see also ref. 18). However, this does not imply that the posteriors lead to predictions that are close to being correct for a given opponent. Our result shows, in fact, that these rational expectations predictions are not close to being correct for almost all opponents.

This still leaves open the possibility that for some combinations of beliefs the players' strategies converge to Nash equilibrium even though their predictions do not. In a repeated game convergence to equilibrium can be given a variety of interpretations; we shall show that the process fails to converge to equilibrium in almost any reasonable sense. Let Q^N be the set of all one-period strategy tuples $q \in \Delta$ such that q occurs in *some* time period in *some* Nash equilibrium of the repeated game. For every $q \in \Delta$ let $d(q, Q^N)$ be the minimum Euclidean distance between q and the compact set Q^N . Given a learning process (f, g) and a specific history h, if the behavioral strategies come close

to Nash equilibrium on h then at a minimum we would expect the following condition to hold,

$$\lim_{T \to \infty} \left[\sum_{t=1}^{T} d(q^{t}(h), Q^{N})^{2} \right] / T = 0.$$
 [2]

This implies that, for every $\varepsilon > 0$, play is within ε of some Nash equilibrium at each point in time except possibly for a sparse set of times. We shall show that the process *fails to come close* to Nash in the sense that Eq. 2 fails to hold for *almost all* histories h.

Theorem. Let v be a continuous density on $[-\lambda/2, \lambda/2]$, and let G be a v-perturbation of a finite, zero-sum, two-person game G^0 , all of whose Nash equilibria have full support. Assume that the players are perfectly rational, have arbitrary discount factors less than unity, and that each updates his predictions of the opponent's future behavior by a learning rule that is based solely on observable actions. If λ is sufficiently small, then for v-almost all payoff realizations, the probability is 1 that someone never learns to predict and that play fails to come close to Nash.

We remark that the set of games for which this impossibility result holds is actually much larger than the one stated in the theorem. Consider, for example, any two-person game G with strategy space $Y_1 \times Y_2$ such that $|Y_i| \ge 2$; all Nash equilibria have full support on $Y_1 \times Y_2$, and every action not in Y_i is dominated strictly by some action in Y_i . Then the theorem holds for perturbed versions of this game. Next let us extend G to an n-person game G^* by adjoining n-2 players as follows: each new player has a strictly dominant action, and G^* is the two-person subgame that results when they play these actions. It follows that for any finite action space $X = \prod X_i$, there exists an n-person game G^* on X such that when the payoffs of G^* are perturbed by small i.i.d. random errors, good prediction fails to occur with probability one.

Now consider any n-person game G on the finite strategy space $X = \Pi X_i$. Suppose that we perturb the payoffs of G by i.i.d. random errors drawn from a normal distribution or in fact any distribution with a continuous density, the support of which is the whole real line. With positive probability the payoffs of the realized game will be close to the game G^* constructed above. Thus as a corollary we obtain the following.

Corollary. Let G be any finite n-person game, the payoffs of which are perturbed once by i.i.d. normally distributed random errors. Assume that the players are perfectly rational, have arbitrary discount factors less than unity, and that each updates his predictions of the opponents' future behavior by a learning rule that is based solely on observable actions. For almost all payoff realizations, there is a positive probability that someone never learns to predict and that play fails to come close to Nash.

Proof of the Theorem. Because the proof is somewhat involved, we shall explain first why the argument given in the introduction for myopic players does not extend easily to the general case. One difficulty is that patient players might interact through conditional strategies that involve no randomization, and these might be predictable at least some of the time. Eliminating this case requires a delicate probabilistic argument. The second difficulty is that even when players randomize and are therefore indifferent among alternative strategies, this does not imply that the stage-game payoffs are solutions of a linear equation. Rather, they are the roots of a nonlinear function, and we must show that the roots of this function constitute a set of measure zero.

To increase the transparency of the proof, we shall give it for

the game of matching pennies. It generalizes readily to any finite zero-sum, two-person game, the stage-game Nash equilibria of which are all strictly interior in the space of mixed strategies. Fix a continuous density ν , the support of which is $[-\lambda/2, \lambda/2]$. To be concrete, we may think of ν as the uniform distribution. The perturbed game has the payoff matrix

where ω_{ij} , ω'_{ij} are i.i.d. random variables distributed according to ν . Fix two rational players, 1 and 2, with discount factors $0 \le \delta_1 \le \delta_2 < 1$. Let their beliefs be f_1 , f_2 , and let their strategies be $g_1(\cdot|A)$, $g_2(\cdot|B)$, where A and B are the realized values of the players' payoff matrices. The functions f_1 , f_2 , g_1 , g_2 will be fixed throughout the proof. All probability statements will be conditional on them without writing this dependence explicitly. Let H(A, B) be the set of all histories h such that good prediction (Eq. 1) holds when the realized payoffs are (A, B). Let P be the set of pairs (A, B) such that good prediction holds with positive probability, that is, $\mu[H(A, B)] > 0$. First we shall show that $\nu(P) = 0$, that is, there are almost no payoff realizations (A, B) such that both players learn to predict with positive probability. In the second part of the proof we shall show that for almost all (A, B) the process fails to come close to Nash.

Lemma 1. For every positive integer m, every $0 < \varepsilon' \le \varepsilon < 1$, and every $(A, B) \in P$, there exists a time T, possibly depending on m, ε , ε' , A, B, such that with μ -probability at least $1 - \varepsilon'$, each player forecasts the other's next-period strategy within ε in each of the periods $T + 1, \ldots, T + m$.

Proof. Let $(A, B) \in P$ and suppose there were no such time T. Then for *every* time T the μ -probability would be greater than $\varepsilon' > 0$ that at least one player misforecasts the opponent's behavior by more than ε in one or more of the periods $T+1,\ldots,T+m$. This would imply that Eq. 1 is violated for almost all histories, that is, $\mu[H(A,B)] = 0$, which contradicts our assumption that $(A,B) \in P$.

Lemma 2. For each $(A, B) \in P$ there exists a time T and a history h^T , possibly depending on A, B, such that conditional on h^T each player's expected future payoffs, discounted to T+1, are bounded above by $c\lambda$ for some positive number c that depends only on the discount rates.

Proof. Given a small $\lambda>0$, choose $m\geq 1$ such that $\delta_2^m\leq \lambda$ and $0<\varepsilon'\leq \varepsilon\leq \lambda/m4^me^m$. As guaranteed by Lemma 1, let h^T be a history such that the μ -probability is at least $1-\varepsilon'$ that each player forecasts the other's next-period strategy within ε in each of the periods $T+1,\ldots,T+m$. Let α_{T+1}^* and β_{T+1}^* be the payoffs that players 1 and 2 *expect* to get from period T+1 on, discounted to period T+1. We shall exhibit a positive constant c, depending only on the discount factors, such that $\alpha_{T+1}^*,\beta_{T+1}^*\leq c\lambda$. Note first that each player has the option of playing 50–50 in each period from T+1 on, which has an expected discounted payoff at least $-\lambda/2$. Because each player's strategy is optimal, it follows that $\alpha_{T+1}^*,\beta_{T+1}^*\geq -\lambda/2$.

For each j, $1 \le j \le m$, let α_{T+j} be player 1's expected undiscounted payoff in period T+j as forecast by player 1 at the end of period T. Define β_{T+j} similarly for player 2. Let H_{j,h^T} be the set of all continuations of h^T to time T+j. Let $\phi_1(h^{T+j})$ denote player 1's probability assessment of $h^{T+j} \in H_{j,h^T}$ and similarly define $\phi_2(h^{T+j})$ for player 2. The true probability is $\mu_0(h^{T+j})$, where μ_0 is μ conditional on h^T . The set of continuations on which someone makes a bad forecast have μ_0 -probability at most ε' . On the remaining good continuations,

loaded at #61431990 on January 7, 2020

each player errs by at most ε in forecasting his opponent's stage-game behavior in each of j stages. Hence for every good continuation h^{T+j} , $|\phi_i(h^{T+j}) - \mu_0(h^{T+j})| \le (1+\varepsilon)^j - 1 \le (j\varepsilon)e^{j\varepsilon}$.

Each player's forecasted payoff in period T+j cannot differ from the actual payoff in period T+j by more than $2+\lambda$ no matter how bad the forecast is. There are 4^j continuations to period T+j including good and bad. Over all of the good ones, player 1's forecasted expected payoff differs from his actual expected payoff by at most $4^j(j\varepsilon)e^{j\varepsilon}$ $(2+\lambda)$. Over all of the bad ones the two differ by at most $\varepsilon'(2+\lambda)$. Thus the difference between 1's forecasted expected payoff, α_{T+j} , and his actual expected payoff, $\bar{\alpha}_{T+j}$, is at most $(\varepsilon'+4^j(j\varepsilon)e^{j\varepsilon})(2+\lambda)$. By assumption, $\varepsilon' \leq \varepsilon \leq \lambda/m4^me^m$ and $j \leq m$, so $\varepsilon'+4^j(j\varepsilon)e^{j\varepsilon} \leq \varepsilon + 4^m(m\varepsilon)e^{m\varepsilon} \leq 2\lambda$. Thus $|\alpha_{T+j}-\bar{\alpha}_{T+j}| \leq 2\lambda(2+\lambda) \leq 6\lambda$. Similarly $|\beta_{T+j}-\bar{\beta}_{T+j}| \leq 6\lambda$. The actual payoffs satisfy $|\bar{\alpha}_{T+j}+\bar{\beta}_{T+j}| \leq \lambda$, from which we conclude that $|\alpha_{T+j}+\beta_{T+j}| \leq 13\lambda$ for $1\leq j \leq m$.

For each j, $1 \le j \le m$, let $\alpha_{T+j}^* = (1 - \delta_1)(\alpha_{T+j} + \delta_1\alpha_{T+j+1} + \delta_1^2\alpha_{T+j+2} + \dots)$ be player 1's expected payoff from period T+j on, discounted to period T+j, as forecast at the end of period T. Similarly define $\beta_{T+j}^* = (1 - \delta_2)(\beta_{T+j} + \delta_2\beta_{T+j+1} + \delta_2^2\beta_{T+j+2} + \dots)$. We claim that α_{T+j}^* , $\beta_{T+j}^* \ge -\lambda/2$ for every j. If not, some player could switch his strategy to a 50–50 random mixture from period T+j on, thus increasing his expected payoff from that time on, which would contradict sequential rationality.

Beyond period T+m, the forecasts may no longer be good within ε . However, neither player expects to get more than $1+\lambda/2$ in any period, thus the sum of expected payoffs beyond period T+m, discounted to period T+1, cannot be more than $(1-\delta_2)\delta_2^m(1+\lambda/2)$. By choice of m, $\delta_2^m \le \lambda$, thus the previous expression is at most 2λ when λ is small. Putting this fact together with $|\alpha_{T+j} + \beta_{T+j}| \le 13\lambda$, it follows that

$$\beta_{T+1}^* \le (1 - \delta_2) \sum_{j=1}^m \delta_2^{j-1} \beta_{T+j} + 2\lambda$$

$$\le (1 - \delta_2) \sum_{j=1}^m \delta_2^{j-1} (13\lambda - \alpha_{T+j}) + 2\lambda$$

$$\le 15\lambda - (1 - \delta_2) \sum_{j=1}^m \delta_2^{j-1} \alpha_{T+j}.$$
[3]

The term $\sum_{j=1}^{m} \delta_{j}^{j-1} \alpha_{T+j}$ is similar in form to α_{T+1}^{*} except that the wrong discount factor is being used, and the sum is truncated. Nevertheless, we claim that if α_{T+1}^{*} is small, then so is the term in question. To see this, consider the identity $\alpha_{T+j}^{*} = \delta_{1}\alpha_{T+j+1}^{*} + (1 - \delta_{1})\alpha_{T+j}$, which holds for all j. From this we obtain

$$\sum_{j=1}^{m} \delta_{2}^{j-1} \alpha_{T+j}^{*} = \delta_{1} \sum_{j=1}^{m} \delta_{2}^{j-1} \alpha_{T+j+1}^{*} + (1-\delta_{1}) \sum_{j=1}^{m} \delta_{2}^{j-1} \alpha_{T+j},$$

and after rearranging terms,

$$\sum_{j=1}^{m} \delta_{2}^{j-1} \alpha_{T+j} = \left[\alpha_{T+1}^{*} + (\delta_{2} - \delta_{1}) \sum_{j=1}^{m-1} \delta_{2}^{j-1} \alpha_{T+j+1}^{*} - \delta_{1} \delta_{2}^{m-1} \alpha_{T+m+1}^{*} \right] / (1 - \delta_{1}).$$
[4]

All of the terms $\alpha_{T+2}^*,\ldots\alpha_{T+m}^*$ are at least $-\lambda/2$, the term α_{T+m+1}^* is at most $1+\lambda/2$, and $\delta_1\delta_2^{m-1} \leq \delta_2^m \leq \lambda$. Thus, the right-hand side of Eq. 4 is bounded below by $\alpha_{T+1}^*/(1-\delta_1)-c'\lambda$, where c'>0 depends only on the discount factors. The left-hand side of Eq. 4 is the summation on the right-hand side

of 3. Substituting this expression into 3 we see that $\beta_{T+1}^* + [(1-\delta_2)/(1-\delta_1)]\alpha_{T+1}^* \le 15\lambda + (1-\delta_2)c'\lambda$. Because $\alpha_{T+1}^*, \beta_{T+1}^* \ge -\lambda/2$, we conclude that both α_{T+1}^* and β_{T+1}^* are bounded above by $c\lambda$ for some c that depends only on the discount factors δ_1 and δ_2 . This concludes the proof of Lemma 2.

Lemma 3. For every positive integer m and all sufficiently small $\lambda > 0$, if $(A, B) \in P$, then there exists a history h^T such that, conditional on h^T at time T, both players randomize in each of the periods $T + 1, \ldots, T + m$.

Proof. As in the proof of Lemma 2, choose $m \ge 1$ such that $\delta_2^m \le \lambda$ and let $0 < \varepsilon \le \lambda/m4^m e^m$. Assume in addition that $\varepsilon' = \varepsilon^{4m}$. Now apply Lemma 1 with 2m instead of m: there is a time T such that, with a probability of at least $1 - \varepsilon'$, the next-period forecasts are within ε of being correct for the periods $T + 1, \ldots, T + 2m$.

For each h^{T+j} , $0 \le j \le 2m-1$, say that h^{T+j} is *good* if both players' next-period forecasts are within ε of being correct; otherwise h^{T+j} is *bad*. Say that h^{T+j} is γ -good if it is good and, conditional on h^{T+j} occurring in period T+j, the probability is at most γ that someone makes a bad next-period forecast in any continuation of h^{T+j} through period T+2m-1.

By choice of T there is at least one state, h^T , that has positive probability under the strategies and is ε' -good. Lemma 2 implies that the expected discounted payoffs from T+1 on are bounded above by $c\lambda$. We claim this implies that both players randomize in period T+1, and in fact each of them chooses each action with probability at least ε . Suppose, to the contrary, that some player (say player 1) chooses action 1 with probability less than ε . Because h^T is good, player 2 forecasts that 1 will play action 2 with probability at least $1-2\varepsilon$. But then player 2 could obtain a higher expected payoff by mismatching (playing action 1) in period T+1 and randomizing fifty-fifty in every period thereafter. (The expected payoff from this strategy is at least $[(1-\delta_2)(1-4\varepsilon)-\lambda/2]$, which is greater than $c\lambda$ for all sufficiently small λ and $\varepsilon \le \lambda$.) This contradiction shows that player 1 chooses each action in period T+1 with probability at least ε , and the same holds for player 2.

It follows that each of the four possible continuations of h^T to period T+1 has probability at least ε^2 . Because $\varepsilon^2 > \varepsilon'$ and h^T is ε' -good, none of these four continuations can be bad, and in fact each of them must be at least $(\varepsilon'/\varepsilon^2)$ -good. Now apply Lemma 2 again (redefining ε' to be $\varepsilon^{4m}/\varepsilon^2$) and conclude that, for every continuation of h^T to some h^{T+2} , the conditional expected payoffs from period T+2 forward are bounded above by $\varepsilon\lambda$. As before, we conclude that both players randomize in period T+2, each putting at least ε on each action. Continuing in this manner, we deduce that both players randomize in every continuation of h^T to period T+2. This concludes the proof of Lemma 3.

The gist of the proof thus far is that, if the payoff realizations (A, B) lead to good predictions with μ -positive probability, then for every sufficiently large positive integer m, there exists a state h^T that induces randomization by both players in each of the next m periods. We now show that this implies that the payoffs are zeroes of a function, the set of zeroes of which has ν -measure zero. This will show that good prediction occurs with ν -measure zero.

Let h^T be any state, and let m be a positive integer. Suppose that player 1 plays action 1 in each of the periods T+1 to T+m, after which he plays an optimal strategy given his beliefs. We can write his expected utility, discounted to time T+1, as a function of his payoff matrix A as follows:

$$U_1(A) = \theta_1 a_{11} + (1 - \delta_1^m - \theta_1) a_{12} + \delta_1^m R_1(A).$$

Here θ_1 comes from player 2's randomization between actions 1 and 2, and the remainder term $R_1(A)$ is convex and bounded. In fact, $|R_1(A)| \le (1 - \delta_1^{\mathrm{m}})(|a_{11}| + |a_{12}| + |a_{21}| + |a_{22}|)$. Similarly define $U_2(A)$ to be player 1's expected utility from playing action

2 for m periods and an optimal strategy thereafter. This can be written analogously to $U_1(A)$ with a remainder function $R_2(A)$ that satisfies the same bound as $R_1(A)$. All of these functions depend of course on h^T .

Ît will be convenient to consider a one-dimensional subspace of the payoff matrices A. Namely, for every four real numbers w, x, y, and z, let $\psi_{x,y,z}(w)$ be the 2×2 matrix with entries $a_{11} = w + x$, $a_{12} = w - x$, $a_{21} = y$, and $a_{22} = z$. Given x, y, and z, define the following function: $F_{x,y,z}(w) = U_1(\psi_{x,y,z}(w)) - U_2(\psi_{x,y,z}(w))$. Abbreviating $\psi_{x,y,z}(w)$ by $\psi(w)$, we can write this in the form

$$F_{x, y, z}(w) = K_{x, y, z} + (1 - \delta_1^m)w + \delta_1^m \{R_1(\psi(w)) - R_2(\psi(w))\},$$

where $K_{x,y,z}$ is a linear function of x,y,z and does not depend on w. The functions $R_i(\psi(w))$ are convex and bounded by the same bound as before. By choosing m to be sufficiently large, we can ensure that $F_{x,y,z}(w)$ is strictly monotone increasing in w. It follows that for any triple x,y,z, there is at most one value of w such that $F_{x,y,z}(w)=0$. Because x,y,z are drawn from the continuous density v, we have $P[\{w:F_{x,y,z}(w)=0|x,y,z\}\}]=0$. By the smoothing theorem (i.e., the law of iterated expectations), it follows that $P[\{(w,x,y,z)\in \mathbf{R}^4:F_{x,y,z}(w)=0\}]=0$.

To state this in terms of the matrix A, let $G(A) = F_{x,y,z}(a_{11} + a_{12})/2)$ where $x = (a_{11} - a_{12})/2$, $y = a_{21}$, and $z = a_{22}$. The preceding implies that $P[\{A: G(A) = 0\}] = 0$. Recalling that F (and thus G) are conditional on a particular history h^T , we can write this as $P[\{A: G(A) = 0\}|h^T] = 0$. Hence, $\sum_{h^T} P[\{A: G(A) = 0\}|h^T] P(h^T) = 0$. In other words, player 1 is only indifferent between actions 1 and 2 on a set of payoff matrices A having ν -measure zero.

Suppose now that (A,B) is a pair for which good prediction holds. Let h^T be a history as guaranteed by Lemma 3, where m is sufficiently large that F is strictly monotone increasing in w. By Lemma 3, player 1 randomizes in each of the periods $T+1,\ldots,T+m$. Hence he is indifferent between playing action 1 or action 2 in each of these periods, an event that has ν -measure zero. We conclude that there are ν -almost no payoff realizations (A,B) such that both players learn to predict with positive probability. This establishes the first claim of the theorem. We also note for future reference that we have actually established the following fact.

Lemma 4. If m is large enough and λ is small enough, then the ν -probability is zero that there exists a state h^T such that, conditional on h^T at time T, both players randomize in each of the periods $T+1,\ldots,T+m$.

It remains to be shown that, for ν -almost all (A, B), play fails to come close to the set of Nash equilibria in the sense that condition 2 fails to hold for almost all histories h. The first step is to show that all Nash equilibria of the repeated game are mixed sufficiently in each time period provided that λ is sufficiently small.

Lemma 5. There exists $\varepsilon > 0$ and $\lambda' > 0$ such that whenever $0 < \lambda \le \lambda'$, every Nash equilibrium of the repeated game puts probability at least 2ε on each action in every time period.

The proof is similar to that of Lemma 2; in outline it runs as follows. In equilibrium, each player's expected discounted payoff

rium, each player's expected discounted payoff supported in part by National Science Foundation Grant SBR 960

- 1. Jordan, J. S. (1991) Games Econ. Behav. 3, 60-91.
- 2. Jordan, J. S. (1993) Games Econ. Behav. 5, 368-386.
- 3. Jordan, J. S. (1995) Games Econ. Behav. 9, 8–20.
- 4. Kalai, E. & Lehrer, E. (1993) Econometrica 61, 1019–1045.
- 5. Lehrer, E. & Smorodinsky, R. (1997) Games Econ. Behav. 18, 116-134.
- 6. Nachbar, J. H. (1997) Econometrica 65, 275–309.
- 7. Nachbar, J. H. (1999) Dept. Economics, Working Paper, Washington Univ., St. Louis.
- 8. Nachbar, J. H. (2001) Soc. Choice Welfare 18, 303-326.
- 9. Binmore, K. (1987) Econ. Philos. 3, 179–214.
- 10. Binmore, K. (1990) Essays on Foundations of Game Theory (Basil Blackwell, Oxford).
- Blume, L. & Easley, D. (1998) in Organizations with Incomplete Information: Essays in Economic Analysis, ed. Majumdar, M. (Cambridge Univ. Press, Cambridge, UK), pp. 61-100

must be at least $-\lambda/2$, because at least this much is guaranteed by randomizing fifty-fifty in every period. Because the actual payoffs in each period sum to λ or less, each player's expected discounted payoff can be bounded from above by $k\lambda$, where k is a positive constant. If some player were to play an action with less than probability 2ε in some period t, the opponent can take a pure action with expected payoff at least $(1 - 4\varepsilon - \lambda/2)$ in period t and get at least $-\lambda/2$ in every period thereafter. When ε and λ are sufficiently small, the expected discounted payoff from such a deviation exceeds $k\lambda$, a contradiction.

Fix $\lambda \in (0, \lambda']$. For each pair of payoff matrices (A, B), let N(A, B) be the set of all histories h such that condition 2 holds, i.e., such that play comes close to Nash in a weak sense. We are going to show that there are μ -almost no such histories for ν -almost all (A, B). This is a consequence of the following.

Lemma 6. Let (A, B) be a pair of payoff matrices such that 2 holds with μ -positive probability. Then for every positive integer m and every sufficiently small λ , there exists a state h^T such that, conditional on h^T , each player randomizes in each of the periods $T + 1, \ldots, T + m$.

By choosing m large enough, it follows from Lemma 4 that there are ν -almost no payoff realizations (A, B) with this property. In other words, for ν -almost all payoff realizations play fails to come close to Nash. Thus, once we establish Lemma 6, we will have completed the proof of the theorem.

Proof of Lemma 6. Fix a pair (A, B) such that condition **2** holds with μ -positive probability. Choose ε and λ such that every element of Q^N puts probability at least 2ε on each action in each time period as guaranteed by Lemma 5. Let m be a positive integer, and let $\varepsilon' = \varepsilon^{2m}$. Let $q^{t+1}(h^t)$ denote the strategies in period t+1 given the history h^t to period t. There exists a time T such that with μ -probability at least $1-\varepsilon'$, $d(q^{t+1}(h^t), Q^N) \le \varepsilon$ for every h^t in the interval $T \le t \le T + m - 1$. (If this were not so, condition **2** would hold with μ -probability zero, contrary to our assumption.)

Say that a history h^{t+1} is *good* if $d(q^{t+1}(h^t), Q^N) \le \varepsilon$. It is *very good* if it is good and all of its successors for the next m periods are good. If a history is good then each action is played in the next period with probability at least ε . Hence every continuation of a good history occurs with probability at least ε^2 . If no history at time T is very good, then the μ -probability of a bad history occurring in the interval $T, T+1, \ldots, T+m-1$ is at least $\varepsilon^{2m-2} > \varepsilon'$, contrary to our assumption. Hence there exists h^T such that $d(q^{t+1}(h^t), Q^N) \le \varepsilon$ for every continuation of h^T in the interval $T \le t \le T + m - 1$, and hence both players randomize for m periods in succession. By Lemma 4 this happens with ν -probability zero. This concludes the proof of Lemma 6, and thereby the proof of the theorem.

We thank the editor, referees, and members of the Santa Fe Institute and the Brookings-Johns Hopkins Center on Social and Economic Dynamics for constructive comments on an earlier draft. This research was supported in part by National Science Foundation Grant SBR 9601743.

- 14. Monderer, D. & Shapley, L. (1996) J. Econ. Theory 68, 258-265.
- 15. Harsanyi, J. (1973) Int. J. Game Theory 2, 1-23.
- 16. Fudenberg, D. & Kreps, D. (1993) Games Econ. Behav. 5, 320-367.
- 17. Nash, J. (1950) Ph.D. Dissertation (Princeton University, Princton).
- 18. Nyarko, Y. (1998) Econ. Theory 11, 643-655.
- Binmore, K. (1993) Studies in Logic and the Foundations of Game Theory: Proceedings of the Ninth International Congress of Logic, Methodology, and the Philosophy of Science, eds. Prawitz, D., Skyrms, B. & Westerstahl, D. (North Holland, Dordrecht), vol. 9, pp. 927-946.

^{12.} Miller, R. & Sanchirico, C. (1997) Columbia University Discussion Paper, No. 9697-25.

Miyasawa, K. (1961) On the Convergence of the Learning Process in a 2 × 2 Non-Zero-Sum Two-Person Game, Economic Research Program, Research Memorandum no. 33, Princeton University, Princeton.