

# Regret Testing: A Simple Payoff-Based Procedure for Learning Nash Equilibrium<sup>1</sup>

DEAN P. FOSTER

Department of Statistics, Wharton School,  
University of Pennsylvania

H. PEYTON YOUNG

Department of Economics, Johns Hopkins University  
and University of Oxford

Running title: Regret Testing

Correspondence: H.P. Young

Department of Economics  
Johns Hopkins University  
3400 N. Charles Street  
Baltimore, MD 21218-2685

Tel: 202 797 6025

pyoung@jhu.edu

This version: February 6, 2006

---

<sup>1</sup>The authors thank Andrew Felton, Sham Kakade, Ben Klemens, Thomas Norman, Kislaya Prasad, and several anonymous referees for constructive comments on an earlier draft. The paper also benefited from audience comments at Universitat Pompeu Fabra, The European University Institute, the University of Siena, and the Second International Congress of the Game Theory Society. An earlier version of this paper appeared as Sante Fe Institute Working Paper 04-12-034.

## Abstract

A learning rule is *uncoupled* if a player does not condition his strategy on the opponent's payoffs. It is *radically uncoupled* if a player does not condition his strategy on the opponent's actions or payoffs. We demonstrate a family of simple, radically uncoupled learning rules whose period-by-period behavior comes arbitrarily close to Nash equilibrium behavior in any finite two-person game.

Keywords: learning, Nash equilibrium, regret, bounded rationality

JEL Classification Numbers: C72, D83.

# 1 Learning equilibrium

A consistent theme of the learning literature is the difficulty of devising rules that converge to Nash equilibrium in general finite games. Hart and Mas-Colell have shown, for example, that there exists no deterministic uncoupled learning process, whose state variable is the joint empirical distribution of play, that converges to Nash equilibrium in every finite game [18]. More recently, they have shown that there exists no *stochastic* uncoupled learning process that is stationary with respect to histories of bounded length, and that guarantees almost sure convergence of the behavioral strategies to Nash equilibrium in every finite game [19].

Implicit in these results is some form of bounded rationality, because fully rational players need not condition their behavior on histories of bounded length. Do analogous impossibility theorems hold when players are Bayesian and fully rational? A well-known result of Kalai and Lehrer would seem to show the opposite: Bayesian rational play converges almost surely to an approximation of Nash equilibrium behavior, provided that players' strategies are absolutely continuous with respect to their opponents' beliefs [22]. Note that Bayesian learning is uncoupled, because a player's posterior beliefs are determined by the opponent's observed pattern of play, not by any direct observation of the opponent's payoffs.

The hitch is that absolute continuity can be very difficult to satisfy when the opponent's payoffs are unknown. The reason, roughly speaking, is that a player's prior must not rule out events that have positive probability under the opponent's strategy, which by assumption must be optimal given the opponent's beliefs. Thus, in effect, absolute continuity requires a form of mutual consistency between the player's prior and the opponent's payoff function [25, 26]. In other words, although Bayesian learning is uncoupled for given beliefs, the absolute continuity condition amounts to an implicit coupling between one person's beliefs and the other person's payoffs. When neither side knows the other's payoff function, this can lead to situations in which the absolute continuity condition fails "generically." In particular,

there exist two-person games of incomplete information such that, given *any* priors over the opponent’s strategy, Bayesian rational play almost surely does not approximate Nash equilibrium behavior, and the absolute continuity condition fails to hold for almost all payoff realizations [11, 20, 21].

Given these negative results, it is not clear whether there exist any plausible learning rules that lead to Nash equilibrium behavior from out-of-equilibrium conditions. In fact, however, the difficulties can be surmounted if one abandons perfect rationality in favor of certain forms of bounded rationality combined with random search. In a previous paper, for example, we showed that Nash equilibrium behavior can be learned by a form of statistical hypothesis testing [12]. In this approach, each player periodically examines the pattern of his opponent’s behavior over recent history, say over the last  $m$  periods, and tests whether his current (probabilistic) model of that behavior is reasonably consistent with actual play. If it is not, he chooses a new model of the opponent from the space of memory- $m$  models, where his choice of new model has a random component that allows the entire model space to be searched with positive probability.

The players are assumed to be boundedly rational in the sense that, at each point in time, they choose smoothed best responses given their current models. Furthermore, they do not update after every play, but only after a sizeable amount of data ( $m$  periods’ worth) has accumulated. With a suitable choice of parameters, it can be shown that the realized behaviors in such a process come close to Nash equilibrium behavior a very large proportion of the time (though they need not *converge* to Nash equilibrium). Notice that the process is not stationary with respect to the length- $m$  histories, because the players’ behavior depends on their current hypotheses, which change over time. Nor is behavior fully rational (though it is certainly purposeful). Thus the learning process dispenses with certain key conditions on which the above-mentioned impossibility results depend.

In this paper we introduce a new type of learning process, called *regret testing*, that also solves the “learning to play Nash” problem but in a different

and even simpler way. Like reinforcement and aspiration models, regret testing depends solely on a player's *realized payoffs* and requires no observation of the opponent or even knowledge of the opponent's existence [2, 4, 5, 7, 8, 23]. It differs from these models by incorporating some degree of random, undirected search. We shall first define the approach informally in order to emphasize its computational simplicity and complete lack of dependence on the actions or payoffs of the opponent. In section 3 we shall define the rule more formally and discuss its connections with other learning rules in greater detail.

## 2 Regret testing

Consider an individual who lives alone. He has  $m$  possible actions, the names of which are written on slips of paper stored in a hat. The hat contains  $h$  papers. Since a given action can be written on multiple papers, the hat is a device for generating probability distributions over actions. Every probability distribution that is expressible in integer multiples of  $1/h$  is represented by one hat. The larger  $h$  is, the more closely any given distribution can be approximated by one of these hats.

*Step 1.* Once each period (say once a minute) he reaches into his current hat, draws a slip, and takes the action prescribed. He then returns the slip to the hat.

*Step 2.* At random times this routine is interrupted by telephone calls. During a call he absent-mindedly chooses an action uniformly at random instead of reaching into the hat.

*Step 3.* Every time he takes an action he receives a payoff. At the end of day  $t$ , he tallies the average payoff,  $\hat{\alpha}_t$ , he received over the course of the day whenever he was not on the phone. For each action  $j$ , he compares  $\hat{\alpha}_t$  with the average payoff,  $\hat{\alpha}_{j,t}$ , he received when he chose  $j$  and was on the phone.

*Step 4.* If at least one of the differences  $\hat{r}_{j,t} = \hat{\alpha}_{j,t} - \hat{\alpha}_t$  is greater than his tolerance level  $\tau > 0$  he chooses a new hat, where each hat has a positive

probability of being chosen. Otherwise he keeps his current hat and the process is repeated on day  $t + 1$ .

Any procedure of this form will be called a *regret testing rule*. The reason is that  $\hat{\alpha}_{j,t}$  amounts to a statistical estimate of the payoff on day  $t$  that the player would have received from playing action  $j$  all day long, hence the difference  $\hat{r}_{j,t} = \hat{\alpha}_{j,t} - \hat{\alpha}_t$  is the *estimated regret* from not having done so.<sup>2</sup> (Recall that the regrets cannot be evaluated directly because the opponent's actions are not observed.) The logic is simple: if one of the payoff-averages  $\hat{\alpha}_{j,t}$  during the experimental periods is significantly larger than the average payoff in the non-experimental periods, the player becomes dissatisfied and chooses a new strategy, i.e., a new hat from the shelf. Otherwise, out of inertia, he sticks with his current strategy.

The revision process (Step 4) allows for many possibilities. The simplest is to choose each hat with *equal* probability, but this lacks behavioral plausibility. Instead, the player could exploit the information contained in the current payoffs, say by favoring strategies (hats) that put high probability on actions with high realized payoff  $\hat{\alpha}_{j,t}$ . Consider, for example, the following revision rule: with probability  $1 - \varepsilon$  adopt the pure strategy that puts probability one on the action  $j$  that maximizes  $\hat{\alpha}_{j,t}$ ; and with probability  $\varepsilon$  choose a strategy at random. This is a trembled form of best response strategy revision, where the tremble is not in the *implementation* of the strategy but in the *choice* of strategy. In particular, a strategy that is far from being a best response strategy can be chosen by mistake, but the probability of such a mistake is small. While the use of recent payoff information may be sensible, however, we do not insist on it. The reason is that the process will eventually approximate Nash equilibrium behavior *irrespective* of the revision rule, as long as every hat is chosen with a probability that is uniformly bounded away from zero at all revision opportunities. This allows for a great deal of latitude in the specification of the learning process.

We hasten to say that this rule is intended to be a contribution to learn-

---

<sup>2</sup>A similar estimation device was used in [9, 16, 17, 18].

ing theory, and should not be interpreted literally as an empirical model of behavior, any more than fictitious play or regret matching should be. Nevertheless it is composed of plausible elements that are found in other learning rules. One key element of regret testing is *inertia*: if there is no particular reason to change, play continues as before. In fact, inertia is built into the rule at two levels: there is no change of strategy while data is being collected over the course of a day, and change is implemented only if a *significant* improvement is possible—in other words, the alternative payoffs must exceed the current average payoff by more than some positive amount  $\tau$ .

Inertia is an important aspect of aspiration learning [2, 4, 7, 23] as well as several other learning rules in the literature, including hypothesis testing [12] and Hart and Mas-Colell’s regret matching [17, 18]. In the latter procedure, a player continues to choose a given action with high probability from one period to the next. When change occurs, the probability of switching to each new action is proportional to its conditional regret relative to the current action.<sup>3</sup> Hart and Mas-Colell show that under this procedure average behavior converges to the set of correlated equilibria (though period-by-period behaviors need not do so).

A second key element of regret testing is that, when a change in strategy occurs, the choice of new strategy has a random component that allows for wide-area *search*. Except for hypothesis testing, this feature is not typical of other learning rules in the literature. For example, under regret matching, a player’s strategy at any given time is either almost pure, or involves switching probabilistically from one almost-pure strategy to another. Similarly, under aspiration learning, a player switches from one pure strategy to an alternative pure strategy when the former fails to deliver payoffs that meet a given aspiration level. In both of these situations there are probabilistic changes among particular classes of strategies, but not a wide-area search among

---

<sup>3</sup>The *conditional regret* of action  $k$  relative to action  $j$  is the increase in average per-period payoff that would have resulted if  $k$  had been played whenever  $j$  actually was played. (The conditional regret is set equal to zero if  $k$  would have resulted in a lower average payoff than  $j$ .)

strategies.

These two elements - inertia and search - play a key role in the learning process. Inertia stabilizes the players' behavior for long enough intervals that the players have a chance to learn something about their opponent's behavior. Search prevents the process from becoming trapped in adjustment cycles, such as the best response cycles that bedevil fictitious play in some settings. Intuitively, the way the process operates is that it discovers a (near) equilibrium through random search, then stays near equilibrium for a long time due to inertia. While it may seem obvious that this ought to work, it is a different matter to show that it actually does work. One difficulty is that the players' search episodes are not independent. Searches are linked via the history of play, so there is no guarantee that the joint strategy space will be searched systematically. A second difficulty is that, even when a search is successful and an equilibrium (or near equilibrium) has been found, the players do not know it. This is because they are ignorant of the opponent's payoff function, hence they cannot tell when a equilibrium is in hand, and may move away again. The essence of the proof is to show that, nevertheless, the expected time it takes to get close to equilibrium is much shorter than the expected time it takes to move away again.

### 3 Formal definitions and main result

Let  $G$  be a two-person game with finite action spaces  $X_1$  and  $X_2$  for players 1 and 2 respectively. Let  $|X_i| = m_i$  and let  $u^i : X_1 \times X_2 \rightarrow \mathbb{R}$  be  $i$ 's utility function. In what follows, we shall always assume (for computational convenience) that the von Neumann Morgenstern utility functions  $u^i(x)$  have been normalized so that all payoffs lie between zero and one:

$$\min_{x \in X_1 \times X_2} u^i(x) \geq 0 \quad \text{and} \quad \max_{x \in X_1 \times X_2} u^i(x) \leq 1. \quad (1)$$

Let  $\Delta_i$  denote the set of probability mixtures over the  $m_i$  actions of player  $i$ . Let  $h_i$  be the size of  $i$ 's hat (a positive integer). The set of distributions in



$\Delta_i$  that are representable as integer multiples of  $1/h_i$  will be denoted by  $P_i$ . Note that every strategy in  $\Delta_i$  can be closely approximated by some strategy in  $P_i$  when  $h_i$  is sufficiently large. Let  $\tau_i > 0$  denote  $i$ 's *tolerance level*, let  $\lambda_i \in (0, 1)$  be the probability that a call is received by a player  $i$  during any given play of the game, and let  $s$  be the number of plays per day.

The *state space* is  $Z = P_1 \times P_2$ , which we shall sometimes refer to as the *probability grid*. The *state* of the learning process at the start of a given day  $t$  is  $z_t = (p_t, q_t) \in P_1 \times P_2$ . For each action  $j$  of player  $i$ , let  $\hat{\alpha}_{j,t}^i = \hat{\alpha}_{j,t}^i(z_t)$  be the average payoff on day  $t$  in those periods when  $i$  played action  $j$  and was on the phone. Let  $\hat{\alpha}_t^i = \hat{\alpha}_t^i(z_t)$  be  $i$ 's average payoff on day  $t$  when *not* on the phone, and let  $\hat{\theta}_t^i = (\hat{\alpha}_t^i, \hat{\alpha}_{1,t}^i, \dots, \hat{\alpha}_{m_i,t}^i)$ . Note that  $\hat{\theta}_t^i$  contains enough information to implement a wide variety of updating rules, including trembled best response behavior, trembled better response behavior, and so forth. Finally, let

$$\hat{r}_t^i(z_t) = \max_{1 \leq j \leq m_i} \hat{\alpha}_{j,t}^i(z_t) - \hat{\alpha}_t^i(z_t). \quad (2)$$

A *regret-testing rule* for player 1 has the following form: there is a number  $\gamma_1 > 0$  such that, for every  $t$ , and every state  $z_t = (p_t, q_t)$ ,

$$\begin{aligned} \hat{r}_t^1(z_t) \leq \tau_1 &\Rightarrow p_{t+1} = p_t \\ \hat{r}_t^1(z_t) > \tau_1 &\Rightarrow P(p_{t+1} = p | p_t, \hat{\theta}_t^1) \geq \gamma_1 \text{ for all } p \in P_1. \end{aligned} \quad (3)$$

The analogous definition holds for player 2. Note that we must have  $\gamma_i \leq 1/|P_i|$  because the conditional probabilities in (3) sum to unity. The case  $\gamma_i = 1/|P_i|$  corresponds to the uniform distribution, that is, all strategies in  $P_i$  are chosen with equal probability when a revision occurs. The class of regret testing rules is more general, however, because it allows for *any* conditional revision probabilities as long as they are uniformly bounded below by some positive constant.

A pair  $(p, q) \in \Delta_1 \times \Delta_2$  is an  $\varepsilon$ -*equilibrium* of  $G$  if neither player can increase his payoff by more than  $\varepsilon$  through a unilateral change of strategy.

**Theorem 1** *Let  $G$  be a finite two-person game played by regret testers and let  $\varepsilon > 0$ . There are bounds on the tolerances  $\tau_i$ , exploration rates  $\lambda_i$ , hat sizes  $h_i$ , and frequency of play  $s$  such that, at all sufficiently large times  $t$ , the players' joint behavior at  $t$  constitutes an  $\varepsilon$ -equilibrium of  $G$  with probability at least  $1 - \varepsilon$ .*

Specific bounds on the parameters are given in Section 5 below.

**Remark 1**

It is not necessary to assume that the players revise their strategies *simultaneously*, that is, at the end of each day. For example, we could assume instead that if player  $i$ 's measured regrets exceed his tolerance  $\tau_i$ , he revises his strategy with probability  $\theta_i \in (0, 1)$ , and with probability  $1 - \theta_i$  he continues to play his current strategy on the following day. One could also assume that the players use *different amounts of information*. Suppose, for example, that player  $i$  looks at the last  $k_i$  days of payoffs ( $k_i$  integer), and revises with probability  $0 < \theta_i < 1$  whenever the estimated regrets exceed  $\tau_i$ . With fixed values of  $k_i$  and  $\theta_i$  this does not change the conclusion of theorem 1 or the structure of the argument in any significant way.

**Remark 2**

Another learning rule with qualitatively similar properties is hypothesis testing [12]. Like regret testing, this approach combines inertia with occasional random search. Unlike regret testing, however, it requires observation of the opponent—indeed the whole object of hypothesis testing is to discern what the opponent's strategy actually is. Regret testing, by contrast, involves no observation of the opponent, no prediction about the opponent, and no optimization against the opponent. It is essentially a *one player rule* that is well-defined whether or not an opponent even exists.

**Remark 3**

Theorem 1 does not assert that the learning process *converges* to an  $\varepsilon$ -

equilibrium of  $G$ ; rather, it says that the players' period-by-period behaviors are *close to equilibrium with high probability* when  $t$  is large. By annealing the learning parameters at a suitable rate, one can achieve convergence in probability to the set of Nash equilibria, as we show in the concluding section. (With some further refinements of the approach one can actually achieve almost sure convergence, as shown in [15]). Although these are probabilistic forms of convergence, the results are quite strong because they hold for the players' period-by-period behaviors. Regret matching, by contrast, only guarantees that the players' *time-average* behaviors converge, and then only to the set of correlated equilibria.<sup>4</sup>

Before giving the proof of theorem 1 in detail, we shall give an overview of some of the technical issues that need to be dealt with. Regret testing defines one-step transition probabilities  $P(z \rightarrow z')$  that lead from any given state  $z$  on day  $t$  to some other state  $z'$  on day  $t + 1$ . Since these transition probabilities do not depend on  $t$ , they define a stationary Markov process  $P$  on the finite state space  $Z$ . A given state  $z = (p, q)$  *induces a Nash equilibrium* in behaviors if and only if the *expected* regrets in that state are non-positive. Similarly,  $(p, q)$  *induces an  $\varepsilon$ -equilibrium* in behaviors if and only if the expected regrets are  $\varepsilon$  or smaller. Note that this is not the same as saying that  $(p, q)$  itself is an  $\varepsilon$ -equilibrium, because the players' behaviors include experimentation, which distorts the probabilities slightly.

If a given state  $z$  does *not* induce an  $\varepsilon$ -equilibrium, the realised regrets  $\widehat{r}_{j,t}^i$  will be larger than  $\varepsilon$  with fairly high probability for at least one of the players. This player will then revise his strategy. Since no strategy on his grid is excluded when he revises, there is a positive probability he will hit upon a strategy that is close to being a best response to the opponent's current strategy. This is not good enough, however, because the new strategy pair does not necessarily induce an  $\varepsilon$ -equilibrium. What must be shown is that

---

<sup>4</sup>Other rules whose long run average behavior converges to the correlated equilibrium set are discussed in [6, 10, 13, 14]. See [27] for a general discussion of the convergence properties of learning rules.

the players arrive *simultaneously* at strategies that induce an  $\varepsilon$ -equilibrium, a point that is not immediately obvious. For example, one player may revise while the second stays put, then the second may revise while the first stays put, and so forth.

Even if they do eventually arrive at an  $\varepsilon$ -equilibrium simultaneously, they must do so in a reasonably short period of time compared to the length of time they stay at the  $\varepsilon$ -equilibrium once they get there. Again this is not obvious. One difficulty is that the players do not know *when* they have arrived—they cannot see the opponent’s strategy, or even his action, so they cannot determine when an  $\varepsilon$ -equilibrium is in hand. In particular, the realized regrets may be large (due to a series of bad draws) even though the state is close to equilibrium (or even *at* an equilibrium), in which case the players will mistakenly move away again. A second difficulty is that revisions by the two players are uncoupled, that is, they cannot coordinate the search process. In reality, however, their searches are linked because the regrets are generated by their joint actions. Thus, the fact that each player conducts a search of his own strategy space whenever he revises need not imply that the *joint* strategy space is searched systematically.

## 4 Entry and Exit Probabilities

The first step in proving theorem 1 is to compare the probability of entering the set of  $\varepsilon$ -equilibrium states with the probability of leaving them. As a preliminary, we need to refine the concept of  $\varepsilon$ -equilibrium as follows. Given a pair of nonnegative real numbers  $(\varepsilon_1, \varepsilon_2)$ , say that a pair of strategies  $(p, q) \in \Delta_1 \times \Delta_2$  is an  $(\varepsilon_1, \varepsilon_2)$ -*equilibrium* if

$$\begin{aligned} \forall p' \in \Delta_1, \quad u^1(p', q) - u^1(p, q) &\leq \varepsilon_1 \\ \forall q' \in \Delta_2, \quad u^2(p, q') - u^2(p, q) &\leq \varepsilon_2. \end{aligned} \tag{4}$$

When  $\varepsilon_1 = \varepsilon_2 = \varepsilon$ , the terms  $\varepsilon$ -equilibrium and  $(\varepsilon_1, \varepsilon_2)$ -equilibrium will be used interchangeably. For any two real numbers  $x, y$  let  $x \wedge y = \min\{x, y\}$  and  $x \vee y = \max\{x, y\}$ .

**Lemma 1** *Let  $m = m_1 \vee m_2, \tau = \tau_1 \wedge \tau_2$  and  $\lambda = \lambda_1 \wedge \lambda_2$ , and suppose that  $0 < \lambda_i \leq \tau/8 \leq 1/8$  for  $i = 1, 2$ . There exist positive constants  $a, b$ , and  $c$  such that, for all  $t$ ,*

- i) If state  $z_t = (p_t, q_t)$  is a  $(\tau_1/2, \tau_2/2)$ -equilibrium, a revision occurs at the end of period  $t$  with probability at most  $ae^{-bs}$  for all  $s$ .*
- ii) If  $z_t$  is not a  $(2\tau_1, 2\tau_2)$ -equilibrium, each player revises at the end of period  $t$  with probability greater than  $1/2$  and both revise with probability greater than  $1/4$ , provided that  $s \geq c$ .*

*It suffices that  $a = 12m$ ,  $b = \lambda\tau^2/256m$ , and  $c = 10^3m^2/\lambda\tau^2$ .*

**Remark:** The proof will show, in addition, that if just *one* of the players, say  $i$ , can increase his payoff by more than  $2\tau_i$ , then  $i$  revises with probability greater than  $1/2$  whenever  $s \geq c$ . Similarly, if one of the players  $i$  cannot increase his payoff by more than  $\tau_i/2$ , then  $i$  revises with probability at most  $ae^{-bs}$ . We shall sometimes use this unilateral version of lemma 1 in what follows.

The proof of lemma 1 involves a straightforward (but somewhat tedious) estimation of tail event probabilities, which is given in the Appendix. While it is a step in the right direction, however, it is not sufficient to establish theorem 1. In particular, it is not enough to know that the process takes a long time (in expectation) to get *out* of a state that is very close to being an equilibrium; we also need to know how long it takes to get *into* such a state from somewhere else. What matters is the *ratio* between these entry and exit probabilities. This issue is addressed by the following general result on stationary, finite Markov chains.

**Lemma 2** *Consider a stationary Markov chain with transition probability function  $P$  on a finite state space  $Z$ . Suppose there exists a nonempty subset of states  $Z^0$  and a state  $w \notin Z^0$  such that:*

- i) *in two periods the process moves from  $w$  into  $Z^0$  with probability at least  $\rho > 0$ ;*
- ii) *once in  $Z^0$  the process stays there for at least one more period with probability at least  $1 - \theta$ .*

*Then for any stationary distribution  $\pi$  of  $P$ ,  $\pi_w \leq 2\theta/\rho$ .*

**Proof:** Let  $\pi$  be a stationary distribution of  $P$ . By definition  $\pi P = \pi$ , hence  $\pi P^2 = \pi$ , that is,  $\pi$  is also a stationary distribution of  $P^2$ . Condition i) of the lemma says that

$$\sum_{z \in Z^0} P^2(w \rightarrow z) \geq \rho. \quad (5)$$

Condition ii) implies that the probability of staying in  $Z^0$  for at least two successive periods is at least  $1 - 2\theta$ , that is,

$$\forall y \in Z^0, \quad \sum_{z \in Z^0} P^2(y \rightarrow z) \geq 1 - 2\theta. \quad (6)$$

Since  $\pi$  is a stationary distribution of  $P^2$ , the stationarity equations imply that

$$\forall z \in Z^0, \quad \sum_{y \in Z^0} \pi_y P^2(y \rightarrow z) + \pi_w P^2(w \rightarrow z) \leq \pi_z. \quad (7)$$

Summing inequality (7) over all  $z \in Z^0$  and using (5) and (6) we obtain

$$(1 - 2\theta) \sum_{y \in Z^0} \pi_y + \pi_w \rho \leq \sum_{z \in Z^0} \pi_z. \quad (8)$$

Hence,

$$\pi_w \rho \leq 2\theta \sum_{z \in Z^0} \pi_z \leq 2\theta. \quad (9)$$

It follows that  $\pi_w \leq 2\theta/\rho$  as claimed.

## 5 Proof of Theorem 1

We begin by restating theorem 1, giving explicit bounds on the parameters. First we need some additional notation. Given  $\delta \geq 0$ , a strategy  $p \in \Delta_1$  is  $\delta$ -dominant for player 1 if

$$\forall p' \in \Delta_1, \forall q \in \Delta_2, u^1(p', q) - u^1(p, q) \leq \delta.$$

The analogous definition holds for player 2. Let  $d(G)$  be the least  $\delta \geq 0$  such that one or both players have a  $\delta$ -dominant strategy. Note that a strategy is 0-dominant if it is a best reply irrespective of the opponent's strategy. (This is slightly weaker than weak dominance, because a 0-dominant strategy is merely as good as any other strategy without necessarily ever being strictly better). Let  $\tau = \tau_1 \wedge \tau_2$ ,  $\lambda = \lambda_1 \wedge \lambda_2$ ,  $\gamma = \gamma_1 \wedge \gamma_2$ , and  $m = m_1 \vee m_2$ .

**Theorem 1 (restatement)** *Let  $G$  be a two-person game on the finite action space  $X = X_1 \times X_2$  and let  $\varepsilon > 0$ . If the players use regret testing with strictly positive parameters satisfying the following bounds, then at all sufficiently large times  $t$  their joint behavior at  $t$  constitutes an  $\varepsilon$ -equilibrium of  $G$  with probability at least  $1 - \varepsilon$ :*

$$\tau_i \leq \varepsilon^2/48 \tag{10}$$

$$\tau_i \leq d^2(G)/48 \quad \text{if } d(G) > 0 \tag{11}$$

$$\lambda_i \leq \tau/16 \tag{12}$$

$$h_i \geq 8\sqrt{m}/\tau \tag{13}$$

$$\gamma_i \leq 1/|P_i(h_i)| \tag{14}$$

$$s \geq (10^3 m^2 / \lambda \tau^2) \ln(10^5 m / \varepsilon^2 \gamma^7). \tag{15}$$

The need for some such bounds may be explained as follows. The tolerances  $\tau_i$  must be sufficiently small relative to  $\varepsilon$  that the players reject with high probability when their behaviors are not an  $\varepsilon$ -equilibrium. The  $\lambda_i$  must be sufficiently small, relative to  $\varepsilon$  and  $\tau$ , that the behaviors are close to equilibrium, and rejection is very unlikely, whenever the state  $(p, q)$  is sufficiently

close to equilibrium. The  $h_i$  must be sufficiently large that the state space actually contains points that are close to equilibrium. The  $\gamma_i$  can be no larger than  $1/|P_i(h_i)|$ , where  $|P_i(h_i)|$  is the number of probability distributions that can be accommodated by a hat of size  $h_i$ . The amount of information collected,  $s$ , must be large enough that the probability of strategy revision is extremely small whenever the behaviors are sufficiently close to equilibrium. In addition,  $s$  must be large enough for Lemma 1 to hold, which is the case under assumption (15). Perhaps the most interesting point, however, is that the tolerances  $\tau_i$  must also be small enough to discriminate in situations where some player has a  $\delta$ -dominant strategy (for small positive  $\delta$ ) but neither player has a 0-dominant strategy, as in (11). The reason is that, when  $d(G) > 0$ , the process may be too sluggish to enter an  $\varepsilon$ -equilibrium and stay there with high probability unless the tolerances are very much smaller than  $d(G)$  in addition to being very much smaller than  $\varepsilon$ .

**Proof of Theorem 1.** In state  $z = (p, q)$ , player 1 is actually playing the strategy  $\tilde{p} = (1 - \lambda_1)p + (\lambda_1/m_1)\vec{1}_{m_1}$ , where  $\vec{1}_{m_1}$  is a length- $m_1$  vector of 1's. Similarly, player 2 is playing  $\tilde{q} = (1 - \lambda_2)q + (\lambda_2/m_2)\vec{1}_{m_2}$ . It follows that if  $(p, q)$  is an  $\varepsilon/2$ -equilibrium of  $G$ , then  $(\tilde{p}, \tilde{q})$  is an  $\varepsilon$ -equilibrium of  $G$  provided that the  $\lambda_i$  are sufficiently small. Since the payoffs lie between zero and one (see assumption (1)), it suffices that  $\lambda_1, \lambda_2 \leq \varepsilon/4$ . This holds because of assumptions (10) and (12).

Let  $E^*$  be the set of states in  $Z$  that actually are  $\varepsilon/2$ -equilibria of  $G$  (ignoring experimentation). We shall show first that, for every stationary distribution  $\pi$  of the process,

$$\sum_{z \notin E^*} \pi_z \leq \varepsilon/2,$$

equivalently,

$$\sum_{z \in E^*} \pi_z \geq 1 - \varepsilon/2. \tag{16}$$

From this and the preceding remark it follows that the players' induced



behaviors  $(\tilde{p}, \tilde{q})$  constitute an  $\varepsilon$ -equilibrium at least  $1 - \varepsilon/2$  of the time (and hence at least  $1 - \varepsilon$  of the time).

We need to show more however: namely, that the behaviors at time  $t$  constitute an  $\varepsilon$ -equilibrium with *probability* at least  $1 - \varepsilon$  for all sufficiently large times  $t$ . To see why this assertion holds, let  $P$  be the transition probability matrix of the process. If the process begins in state  $z_0$ , then the probability of being in state  $z$  at time  $t$  is  $P^t(z_0 \rightarrow z)$ , where  $P^t$  is the  $t$ -fold product of  $P$ . We claim that  $P$  is acyclic; indeed this follows from the fact that for any state  $z$ ,  $P(z \rightarrow z) > 0$ . (Recall that, whenever a player revises, he chooses his previous strategy with positive probability.) It follows from standard results that the following limit exists

$$\forall z \in Z, \quad \lim_{t \rightarrow \infty} P^t(z_0 \rightarrow z) = \pi_z, \quad (17)$$

and the limiting distribution  $\pi$  is a stationary distribution of  $P$  [24, Theorem 1.2]. From this and (16) it follows that

$$\lim_{t \rightarrow \infty} \sum_{z \notin E^*} P^t(z_0 \rightarrow z) \leq \varepsilon/2. \quad (18)$$

Hence

$$\exists T \forall t \geq T, \quad \sum_{z \in E^*} P^t(z_0 \rightarrow z) \geq 1 - \varepsilon. \quad (19)$$

Thus, for all  $t \geq T$ , the probability is at least  $1 - \varepsilon$  that  $z_t \in E^*$ , in which case the induced behaviors at time  $t$  form an  $\varepsilon$ -equilibrium of  $G$ . This is precisely the desired conclusion. It therefore suffices to establish (15) to complete the proof of theorem 1. We shall consider two cases:  $d(G) > 0$  and  $d(G) = 0$ .

**Case 1**  $d(G) > 0$ : neither player has a 0-dominant strategy.

For every pair  $(p, q) \in \Delta_1 \times \Delta_2$ , there exists  $(p', q') \in Z$  such that

$$|p' - p| \leq \sqrt{m_1}/h_1 \text{ and } |q' - q| \leq \sqrt{m_2}/h_2. \quad (20)$$

By the lower bound (13) on the  $h_i$ , it follows that there is a point  $(p', q') \in Z$  such that

$$|p' - p| \leq \tau_1/8 \text{ and } |q' - q| \leq \tau_2/8. \quad (21)$$

Now let  $(p, q)$  be a Nash equilibrium in the full space of mixed strategies,  $\Delta_1 \times \Delta_2$ . By (21) there is a state  $e^* = (p^*, q^*) \in Z$  such that  $|p^* - p| \leq \tau_1/8$  and  $|q^* - q| \leq \tau_2/8$ . Since all payoffs are bounded between zero and one,  $e^*$  is a  $(\tau_1/8, \tau_2/8)$ -equilibrium. In particular,  $e^* \in E^*$ , because by (10),  $\tau_1/8, \tau_2/8 \leq \varepsilon/2$ . We shall fix  $e^* = (p^*, q^*)$  for the remainder of the proof of case 1.

It follows from Lemma 1, part (i), that

$$P(e^* \rightarrow e^*) \geq 1 - ae^{-bs}. \quad (22)$$

The next step is to show that for all  $w \notin E^*$ , the process enters  $E^*$  in two periods with fairly high probability; then we shall apply Lemma 2.

**Case 1a.**  $w \notin E^*$  and each player can, by a unilateral deviation, increase his payoff by more than  $\varepsilon/2$ .

Suppose that  $z_t = w = (p, q)$ . Since each player  $i$  can increase his payoff by more than  $\varepsilon/2$ , he can certainly increase it by more than  $2\tau_i$  (because of the bound  $\tau_i \leq \varepsilon^2/48$ ). It follows from Lemma 1, part (ii) that the probability is at least  $1/4$  that both players revise at the end of day  $t$ .

Conditional on both revising, the probability is at least  $\gamma^2$  that player 1 chooses  $p^*$  and player 2 chooses  $q^*$  in period  $t + 1$ . Hence

$$P(w \rightarrow e^*) \geq \gamma^2/4, \quad (23)$$

so by (22),

$$P^2(w \rightarrow e^*) \geq (\gamma^2/4)(1 - ae^{-bs}). \quad (24)$$

**Case 1b.**  $w \notin E^*$  and only one of the players can improve his payoff by more than  $\varepsilon/2$ .

This case requires a two-step argument: we shall show that the process can transit from state  $w$  to some intermediate state  $x$  with the property that each player  $i$  can increase his payoff by more than  $2\tau_i$ . As in the proof of Case 1a, it follows that  $P(x \rightarrow e^*) \geq \gamma^2/4$ .

We now establish the existence of such an intermediate state. Assume without loss of generality that in state  $w = (p, q)$ , player 1 can increase his payoff by more than  $\varepsilon/2$ , whereas player 2 cannot. In particular, if  $p' \in \Delta_1$  is a best response to  $q$ , then

$$u^1(p', q) - u^1(p, q) > \varepsilon/2. \quad (25)$$

Let  $\delta = d(G)$ : by definition neither player has a  $\delta'$ -dominant strategy for any  $\delta' < \delta$ . In particular,  $q$  is not  $\delta/2$ -dominant for player 2. Hence there exists  $p^0 \in \Delta_1$  and  $q' \in \Delta_2$  such that

$$u^2(p^0, q') - u^2(p^0, q) > \delta/2. \quad (26)$$

Consider the strategy

$$p'' = (\delta/4)p + (1 - \delta/4)p^0. \quad (27)$$

By assumption,  $p'$  is a best response to  $q$ , so  $u^1(p', q) - u^1(p^0, q) \geq 0$ . It follows from (25) and (27) that

$$\begin{aligned} u^1(p', q) - u^1(p'', q) &= (\delta/4)[u^1(p', q) - u^1(p, q)] \\ &\quad + (1 - \delta/4)[u^1(p', q) - u^1(p^0, q)] \\ &\geq (\delta/4)[u^1(p', q) - u^1(p, q)] \\ &> \delta\varepsilon/8. \end{aligned} \quad (28)$$

By assumptions (10) and (11),  $\tau_1 \leq \delta^2/48$  and  $\tau_1 \leq \varepsilon^2/48$ , hence  $48\tau_1 \leq \delta\varepsilon$ , which implies  $6\tau_1 < \delta\varepsilon/8$ . From this and (28) we conclude that, given  $(p'', q)$ , player 1 can deviate and increase his payoff by more than  $6\tau_1$ .

For player 2 we have, by definition of  $p''$ ,

$$\begin{aligned} u^2(p'', q') - u^2(p'', q) &= (\delta/4)[u^2(p, q') - u^2(p, q)] \\ &\quad + (1 - \delta/4)[u^2(p^0, q') - u^2(p, q)]. \end{aligned}$$

Since utilities are bounded between 0 and 1, the first term on the right-hand side is at least  $-\delta/4$ . The second term is greater than  $(1 - \delta/4)(\delta/2) > 3\delta/8$ , by (26). Hence

$$u^2(p'', q') - u^2(p'', q) > \delta/8. \quad (29)$$

Since  $\tau_2 \leq \delta^2/48 < \delta/48$ , player 2 can deviate from  $(p'', q)$  and increase his payoff by more than  $6\tau_2$ . Hence  $(p'', q)$  is not a  $(6\tau_1, 6\tau_2)$ -equilibrium.

Although  $q$  is on player 2's grid, the definition of  $p''$  in (27) does not guarantee that it is on player 1's grid. We know, however, that there exists a grid point  $(p''', q)$  such that  $|p''' - p''| \leq \sqrt{m_1}/h_1$ . Since all payoffs lie between zero and one, the difference in payoff between  $(p''', q)$  and  $(p'', q)$  is at most  $\sqrt{m_1}/h_1$  for *both* players. From (13) it follows that  $\sqrt{m_1}/h_1 \leq \tau/8 \leq \tau_i/8$  for both players ( $i = 1, 2$ ). Since  $(p'', q)$  is not a  $(6\tau_1, 6\tau_2)$ -equilibrium, it follows that  $(p''', q)$  is not a  $(5\tau_1, 5\tau_2)$ -equilibrium (and is on the grid).

Let  $x = (p''', q)$ . As in the proof of Case 1a, it follows that  $P(x \rightarrow e^*) \geq \gamma^2/4$ . Further, the process moves from state  $w$  to state  $x$  with probability at least  $\gamma/2$ , because only player 1 needs to revise:  $w$  and  $x$  differ only in the first coordinate. Hence,

$$P^2(w \rightarrow e^*) \geq \gamma^3/8. \quad (30)$$

In case 1a we found that  $P^2(w \rightarrow e^*) \geq (\gamma^2/4)(1 - ae^{-bs})$ , which is at least  $\gamma^2/8$  provided that  $ae^{-bs} \leq 1/2$ . This certainly holds under the assumptions in Lemma 1 on  $a$ ,  $b$ , and  $s$ . Since  $\gamma^2/8 \geq \gamma^3/8$ , it follows that in *both* cases

$$\forall w \notin E^*, \quad P^2(w \rightarrow e^*) \geq \gamma^3/8. \quad (31)$$

In both cases we also have  $P(e^* \rightarrow e^*) \geq 1 - ae^{-bs}$ , by (22). Now apply Lemma 2 with  $Z^0 = \{e^*\}$ ,  $\rho = \gamma^3/8$ , and  $\theta = ae^{-bs}$ . We conclude that in *both*

case 1a and case 1b, for every stationary distribution  $\pi$  of  $P$ ,

$$\forall w \notin E^*, \quad \pi_w \leq \frac{2ae^{-bs}}{\gamma^3/8} = 16ae^{-bs}/\gamma^3. \quad (32)$$

There are at most  $1/\gamma^2$  states in  $Z$  altogether, so

$$\sum_{w \notin E^*} \pi_w \leq 16ae^{-bs}/\gamma^5. \quad (33)$$

The right-hand side will be at most  $\varepsilon/2$  if  $ae^{-bs} \leq \gamma^5\varepsilon/32$ , that is, if

$$s \geq (1/b) \ln(32a/\gamma^5\varepsilon). \quad (34)$$

By Lemma 1, we can take  $a = 12m$  and  $b = \lambda\tau^2/256m$ . Thus it suffices that

$$s \geq \frac{256m}{\lambda\tau^2} \ln(384m/\gamma^5\varepsilon),$$

which is implied by the stronger bound in (15). This concludes the proof of Case 1.

**Case 2.**  $d(G) = 0$ ; some player has a 0-dominant strategy.

Fix a probability  $0 < \beta < \frac{1}{2}$  that is *much smaller* than  $\gamma$  and *much larger* than  $ae^{-bs}$ ; later we shall specify  $\beta$  and  $s$  more exactly. Define the following subset of states:

$$Z^\beta = \{z = (p, q) : \forall t, \forall q' \in P_2, \quad P(p_{t+1} \neq p \mid z_t = (p, q')) \leq \beta\}.$$

In words,  $Z^\beta$  is the set of states such that the first player changes strategy with probability at most  $\beta$  no matter what strategy the second player is using on his grid.

Without loss of generality assume that player 1 has a 0-dominant strategy. Then he has a *pure* 0-dominant strategy, say  $p^*$ , which is in  $P_1$ . We shall fix  $p^*$  for the remainder of the proof.

Let  $Z^*$  be the set of states whose first coordinate is  $p^*$ . Then player 1 rejects with probability at most  $ae^{-bs}$  (see the remark after lemma 1), that is,

$$P(z_{t+1} \in Z^* \mid z_t \in Z^*) \geq 1 - ae^{-bs}. \quad (35)$$

Hence  $Z^* \subseteq Z^\beta$  provided that  $ae^{-bs} \leq \beta$ , which holds whenever  $s$  is sufficiently large (we shall assume henceforth that this is the case).

Let  $w \in Z - E^*$ . There are two possibilities:  $w \notin Z^\beta$  and  $w \in Z^\beta$ .

**Case 2a.**  $w \notin E^*$  and  $w \notin Z^\beta$ .

Since  $w = (p, q) \notin E^*$ ,  $w$  is not an  $\varepsilon/2$ -equilibrium, and hence is not a  $(2\tau_1, 2\tau_2)$ -equilibrium. By lemma 1 the probability is at least  $1/2$  that there will be a revision next period by at least one of the players. If player 1 revises, a transition of form  $(p, q) \rightarrow (p^*, \cdot) \in Z^*$  occurs with probability at least  $\gamma$ . After that,  $(p^*, \cdot)$  stays in  $Z^*$  for one more period with probability at least  $1 - ae^{-bs}$ , which is at least  $1/2$  because  $ae^{-bs} < \beta < 1/2$ . Hence in this case

$$P^2(w \rightarrow Z^*) \geq \gamma/4. \quad (36)$$

If player 1 does not revise but player 2 does, then with probability at least  $\gamma$  we have a transition of form  $(p, q) \rightarrow (p, q')$ , where  $q' \in P_2$  is a strategy for player 2 that will make player 1 revise with probability greater than  $\beta$ . (There is such a  $q'$  because of our assumption that  $w \notin Z^\beta$ .) In the following period the transition  $(p, q') \rightarrow (p^*, \cdot)$  occurs with probability greater than  $\beta\gamma$ . Hence in this case

$$P^2(w \rightarrow Z^*) \geq \beta\gamma^2/2. \quad (37)$$

Therefore, in either case,

$$P^2(w \rightarrow Z^*) \geq (\beta\gamma^2/2 \wedge \gamma/4) \geq \beta\gamma^2/4. \quad (38)$$

Now apply lemma 2 with  $Z^0 = Z^*$ ,  $\theta = ae^{-bs}$ , and  $\rho = \beta\gamma^2/4$ . Since  $w \notin Z^*$  we conclude that

$$\pi_w \leq 2ae^{-bs}/(\beta\gamma^2/4) = 8ae^{-bs}/\beta\gamma^2. \quad (39)$$

**Case 2b.**  $w \notin E^*$  and  $w \in Z^\beta$ .

By definition of  $Z^\beta$ , player 1 revises with probability at most  $\beta$ , which by assumption is less than  $1/2$ . Since  $w = (p, q) \notin E^*$ , some player  $i$  can

increase his payoff by at least  $\varepsilon/2$ , and hence by more than  $2\tau_i$ . This player will revise with probability greater than  $1/2$  (see the remark after lemma 1), hence  $i$  cannot be player 1 (who revises with probability less than  $1/2$ ). Therefore  $i$  must be player 2. By (13), there exists  $q''$  on player 2's grid that is within  $\tau_2/8$  of a best response to  $p$ . The probability is at least  $\gamma$  that 2 chooses  $q''$  when he revises. Putting all of this together, we conclude that

$$P((p, q) \rightarrow (p, q'')) \geq \gamma/4. \quad (40)$$

By construction, state  $(p, q'')$  is a  $(\cdot, \tau_2/8)$ -equilibrium, hence player 2 revises with probability at most  $ae^{-bs}$  (see the remark after lemma 1). By assumption,  $(p, q) \in Z^\beta$ , so player 1 revises with probability at most  $\beta$  against any strategy of player 2, including  $q''$ . Hence  $(p, q'')$  is also in  $Z^\beta$ , and

$$P((p, q'') \rightarrow (p, q'')) \geq (1 - \beta)(1 - ae^{-bs}) \geq (1 - \beta)^2 > 1 - 2\beta. \quad (41)$$

From this and (40) we have

$$P^2((p, q) \rightarrow (p, q'')) \geq (\gamma/4)(1 - \beta)^2 > \gamma/16,$$

the latter since  $\beta < \frac{1}{2}$ . Now apply lemma 2 with  $Z^0 = \{(p, q'')\}$ ,  $\rho = \gamma/16$ , and  $\theta = 2\beta$ . It follows that for every stationary distribution  $\pi$  of  $P$ ,

$$\pi_w \leq 2(2\beta)/(\gamma/16) = 64\beta/\gamma. \quad (42)$$

Combining (39) and (42), it follows that in both case 2a and case 2b,

$$\forall w \notin E^*, \quad \pi_w \leq 64\beta/\gamma \vee 8ae^{-bs}/\beta\gamma^2. \quad (43)$$

The size of the state space is at least  $1/\gamma^2$ . Summing (43) over all  $w \notin E^*$  it follows that

$$\pi(Z - E^*) \leq (1/\gamma^2)(64\beta/\gamma \vee 8ae^{-bs}/\beta\gamma^2). \quad (44)$$

We wish to show that this is at most  $\varepsilon/2$ . This will follow if we choose  $\beta$  and  $s$  so that  $64\beta/\gamma^3 = \varepsilon/4$  and  $8ae^{-bs}/\beta\gamma^4 \leq \varepsilon/4$ . Specifically, it suffices that

$$\beta = \varepsilon\gamma^3/256 \quad (45)$$

and

$$s \geq (1/b) \ln(8192a/\varepsilon^2\gamma^7). \quad (46)$$

By Lemma 1 we may choose  $a = 12m$  and  $b = \lambda\tau^2/256m$ , hence it suffices that,

$$s \geq (256m/\lambda\tau^2) \ln(98,304m/\varepsilon^2\gamma^7).$$

This certainly holds under (15), which states that  $s \geq (10^3m^2/\lambda\tau^2) \ln(10^5m/\varepsilon^2\gamma^7)$ . This concludes the proof of the theorem.

## 6 Convergence in probability

Theorem 1 says that, for a given game  $G$ , regret testing induces an  $\varepsilon$ -equilibrium with high probability provided that the learning parameters satisfy the bounds given in (10)-(15). But it does not imply that, for a given set of parameters, an  $\varepsilon$ -equilibrium occurs with high probability for all games  $G$ . The difficulty is condition (11), which in effect requires that  $d(G)$  not fall into the interval  $(0, \sqrt{48(\tau_1 \vee \tau_2)})$ . If we think of  $G$  as a vector of  $2m_1m_2$  payoffs in Euclidean space, the excluded set will be small relative to Lebesgue measure whenever the  $\tau_i$  are small. Thus, if we tighten  $\tau_1, \tau_2$  and the other parameters in tandem, the learning process will eventually capture all games in the “net,” that is, there will be no excluded cases. In this section we shall show even more, namely, that by tightening the parameters sufficiently slowly, the players’ period-by-period behavioral strategies converge in probability to the set of Nash equilibria of  $G$ .

Fix an  $m_1 \times m_2$  action space  $X = X_1 \times X_2$  and consider all games  $G$  on  $X$  with payoffs normalized to lie between zero and one. For each  $\varepsilon > 0$ , we shall choose particular values of the parameters that satisfy all the bounds except (11), namely,

$$\tau_i(\varepsilon) = \varepsilon^2/48, \quad (47)$$



$$\lambda_i(\varepsilon) = \tau/16 \quad (48)$$

$$h_i(\varepsilon) = \lceil 8\sqrt{m}/\tau \rceil, \quad (49)$$

$$\gamma_i(\varepsilon) = 1/|P_i(h_i(\varepsilon))| \quad (50)$$

$$s(\varepsilon) = \lceil (10^3 m^2 / \lambda \tau^2) \ln(10^5 m / \varepsilon^2 \gamma^7) \rceil. \quad (51)$$

Recall that  $|P_i(h_i(\varepsilon))|$  is the number of distributions on  $i$ 's grid when his hat size is  $h_i(\varepsilon)$ . Hence (50) implies that each player chooses a new hat with *uniform* probability whenever a revision is called for. This will prove to be analytically convenient in what follows, although more general assumptions could be made.

Let  $P_G(\varepsilon)$  denote the finite-state Markov process determined by  $G$  and the parameters  $(\tau_1(\varepsilon), \dots, s(\varepsilon))$ . Let  $\mathcal{E}_G(\varepsilon)$  be the finite subset of states that induce an  $\varepsilon$ -equilibrium of  $G$ .

**Definition 1** *Let  $P$  be an acyclic, finite Markov process and  $\mathcal{A}$  a subset of states. For each  $\varepsilon > 0$ , let  $T(P, \mathcal{A}, \varepsilon)$  be the first time (if any) such that, for all  $t \geq T(P, \mathcal{A}, \varepsilon)$  and all initial states, the probability is at least  $1 - \varepsilon$  that the process is in  $\mathcal{A}$  at time  $t$ .*

It follows from theorem 1 that  $T(P_G(\varepsilon), \mathcal{E}_G(\varepsilon), \varepsilon)$  is finite for all games  $G$  such that  $d(G) \notin (0, \sqrt{48(\tau_1 \vee \tau_2)})$ . By assumption (47), this holds whenever  $d(G) \notin (0, \varepsilon)$ . In this case, for all  $t \geq T(P_G(\varepsilon), \mathcal{E}_G(\varepsilon), \varepsilon)$ , the probability is at least  $1 - \varepsilon$  that the behavioral strategies constitute an  $\varepsilon$ -equilibrium of  $G$  at time  $t$ .

The time  $T(P_G(\varepsilon), \mathcal{E}_G(\varepsilon), \varepsilon)$  may depend on the payoffs, because these affect the details of the transition probabilities and the states that correspond to  $\varepsilon$ -equilibria of  $G$ . We claim, however, that for every  $\varepsilon > 0$  there is a time  $T(\varepsilon)$  such that  $T(\varepsilon) \geq T(P_G(\varepsilon), \mathcal{E}_G(\varepsilon), \varepsilon)$  for all  $G$  such that  $d(G) \notin (0, \varepsilon)$ .

To see why this is so, consider the realization of plays on any given day. A realization is a sequence of  $s(\varepsilon)$  actions for each player and a sequence of  $s(\varepsilon)$  binary outcomes (say 0 or 1) that indicate whether a given action was taken

by that player while on the phone or not. Hence there are  $(4m_1m_2)^{s(\varepsilon)}$  possible realizations. We may partition them into four disjoint classes: sequences that are rejected by both players, sequences that are rejected by player 1 but not player 2, sequences that are rejected by player 2 but not player 1, and sequences that are accepted by both. (Notice that this partition does not depend on the day  $t$  or on the strategies  $(p_t, q_t)$  in force during that day, but it does depend on the game  $G$ .) However, a player's response given a rejection does not depend on the sequence because we are assuming that each player chooses a new strategy with uniform probability over all distributions on his grid.

The number of length- $s(\varepsilon)$  realizations is finite, and there are finitely many ways of partitioning them into four classes. Further, the probability that each sequence will be realized on a given day  $t$  is determined by the state  $(p_t, q_t)$ , and there are finitely many states. Hence, over all  $G$ , there can be only a finite number of Markov transition matrices  $P_G(\varepsilon)$ . Further, there are finitely many subsets of states that can be used to define  $\mathcal{E}_G(\varepsilon)$ . Let us enumerate all possible pairs  $(P_G(\varepsilon), \mathcal{E}_G(\varepsilon))$  as follows  $(P_1, \mathcal{E}_1), \dots, (P_k, \mathcal{E}_k)$ . Now define  $T(\varepsilon) = \max_{1 \leq j \leq k} T(P_j, \mathcal{E}_j, \varepsilon)$ . Then  $T(\varepsilon)$  has the property that, for all  $G$  such that  $d(G) \notin (0, \varepsilon)$ , and for all  $t \geq T(\varepsilon)$ , the behavioral strategies constitute an  $\varepsilon$ -equilibrium at time  $t$  with probability at least  $1 - \varepsilon$ .

**Definition 2 (annealed regret testing)** *Consider any positive sequence  $\varepsilon_1 > \varepsilon_2 > \varepsilon_3 > \dots$  decreasing to zero. The annealed regret testing procedure at stage  $k$  is the regret testing procedure with parameters  $(\tau_1(\varepsilon_k), \tau_2(\varepsilon_k), \lambda_1(\varepsilon_k), \lambda_2(\varepsilon_k), h_1(\varepsilon_k), h_2(\varepsilon_k), \gamma_1(\varepsilon_k), \gamma_2(\varepsilon_k), s(\varepsilon_k))$  as in (47)-(51). Each day that the process is in stage  $k$ , the probability of moving to stage  $k+1$  on the following day is*

$$p_k \equiv \frac{\varepsilon_{k+1}^2}{2k^2 T(\varepsilon_{k+1})} \quad (52)$$

**Theorem 2** *Fix an  $m_1 \times m_2$  action space  $X = X_1 \times X_2$ . Annealed regret testing has the property that, for every game  $G$  on  $X$ , the behavioral strategies converge in probability to the set of Nash equilibria of  $G$ .*

Although annealed regret testing seems to require that each player has an arbitrarily long memory, the process is actually of much lower dimension. To see why, let us fix a particular player  $i$ . Create one “payoff register” and one “counting register” for each of  $i$ ’s actions, plus one general payoff register and one general counting register for all of his actions together. Let  $k$  be a state variable that is common to all the players and indicates what stage the process is in (i.e., what parameters are currently in force). At each time  $t$ ,  $i$ ’s general payoff register contains the running total of the payoffs he received when not on the phone, and the general counting register contains the number of times he was not on the phone. Similarly, each action-specific register contains the running total of the payoffs he received when he was on the phone and played that action, and the number of times the action was played while on the phone. When the sum over all counting registers reaches  $s(\varepsilon_k)$ , player  $i$  conducts a test using the  $k^{\text{th}}$  set of parameters, revises his strategy if this is called for, and empties all the registers. The process is then repeated. Thus player  $i$  needs only to keep track of  $2m_i + 3$  numbers – two for each of his actions, two in the general registers, and the current stage  $k$ . Thus the learning process requires very little memory or computational sophistication.

**Proof of theorem 2:** Given  $G$ , it suffices to show that, for every  $\varepsilon > 0$ , there is a finite time  $T_\varepsilon$  (possibly depending on  $G$ ) such that for all  $t \geq T_\varepsilon$ , the probability is at least  $1 - \varepsilon$  that the behavioral strategies  $(\tilde{p}_t, \tilde{q}_t)$  constitute an  $\varepsilon$ -equilibrium of  $G$ . Indeed, for every  $\delta > 0$  there exists  $0 < \varepsilon_\delta \leq \delta$  such that every  $\varepsilon_\delta$ -equilibrium lies within  $\delta$  of the compact set  $\mathcal{N}_G$  of Nash equilibria of  $G$ . Hence, for all  $t \geq T_{\varepsilon_\delta}$ , the probability is at least  $1 - \varepsilon_\delta \geq 1 - \delta$  that  $(\tilde{p}_t, \tilde{q}_t)$  lies within  $\delta$  of  $\mathcal{N}_G$ . Thus, the behavioral strategies converge in probability to the set of Nash equilibria of  $G$ .

To facilitate the proof we will define three integer-valued random variables  $N_t$ ,  $T_k$ , and  $W_k$  that describe the process as it transitions through stages. Let  $N_t$  be the stage that the process is in on day  $t$ . In other words, on day  $t$  the

process is using the parameters  $(\tau_1(\varepsilon_{N_t}), \tau_2(\varepsilon_{N_t}), \lambda_1(\varepsilon_{N_t}), \lambda_2(\varepsilon_{N_t}), h_1(\varepsilon_{N_t}), h_2(\varepsilon_{N_t}), \gamma_1(\varepsilon_{N_t}), \gamma_2(\varepsilon_{N_t}), s(\varepsilon_{N_t}))$ . The distribution of the realizations of  $N_t$  depends on the transition probabilities as follows:

$$\begin{aligned} N_1 &= 1 \\ N_{t+1} &= \begin{cases} N_t & \text{with probability } 1 - p_{N_t} \\ N_t + 1 & \text{with probability } p_{N_t} \end{cases} \end{aligned}$$

The first time that the system uses the  $k^{\text{th}}$  set of parameters will be denoted by  $T_k$ , that is,  $T_k \equiv \inf_t \{t : N_t \geq k\}$ . Now define  $W_t \equiv t - T_{N_t}$  to be the length of time since the parameters were last changed. In essence, the proof consists of establishing two facts about  $W_t$ . First we will show that if  $W_t$  is “large” for a given  $t$ , the behavioral strategies are nearly a Nash equilibrium with high probability at time  $t$ . This follows by applying theorem 1 to this setting. Second, we will show that the probability that  $W_t$  is “large” converges to one as  $t$  converges to infinity. This follows from our assumption that the transition probabilities  $p_k$  are small. We now establish these points in detail.

For any game  $G$  on  $X$ , if  $d(G) > 0$  then  $d(G) \geq \varepsilon_k$  for all sufficiently large  $k$ , because the sequence  $\{\varepsilon_k\}$  decreases to zero. The least such  $k$  will be called the *critical index* of  $G$ , and denoted by  $k_G$ . In case  $d(G) = 0$ , we will take  $k_G = 1$ . Fix  $\varepsilon > 0$ . Define  $k_G^* = k_G \vee \min_k \{k \mid \varepsilon_k \leq \varepsilon/4\}$ . It follows that if  $N_t \geq k_G^*$  then  $\varepsilon_{N_t} \leq \varepsilon/4$  and  $d(G) \geq \varepsilon_{N_t}$ .

Since  $N_t \rightarrow \infty$  almost surely as  $t \rightarrow \infty$ , there is a time  $T^*$  such that, for all  $t \geq T^*$ , the probability is at least  $1 - \varepsilon/4$  that  $N_t \geq k_G^*$ . From now on we shall only consider  $t \geq T^*$ .

Given  $t \geq T^*$ , consider two cases:  $W_t \geq T(\varepsilon_{N_t})$  and  $W_t < T(\varepsilon_{N_t})$ . In the first case, the process is an  $\varepsilon_{N_t}$ -equilibrium with probability at least  $1 - \varepsilon_{N_t}$ . Since  $t \geq T^*$ ,  $N_t \geq k_G^*$  with probability at least  $1 - \varepsilon/4$ , in which case  $\varepsilon_{N_t} \leq \varepsilon/4$ . It follows that at time  $t$  the process is in an  $\varepsilon/4$ -equilibrium, and hence an  $\varepsilon$ -equilibrium, with probability at least  $(1 - \varepsilon/4)(1 - \varepsilon/4) \geq 1 - \varepsilon/2$ .

To complete the proof, it therefore suffices to show that, for all sufficiently large  $t$ , the second case occurs with probability at most  $\varepsilon/2$ , that is, there exists  $T^{**}$  such that

$$\forall t \geq T^{**}, \quad P(W_t < T(\varepsilon_{N_t})) \leq \varepsilon/2. \quad (53)$$

To establish (53) we proceed as follows. Recall that in the  $k^{\text{th}}$  stage of the process, the parameter values are  $(\tau_1(\varepsilon_k), \dots, s(\varepsilon_k))$ . By choice of  $p_k$ , the  $k^{\text{th}}$  stage lasts for  $2k^2T(\varepsilon_{k+1})/\varepsilon_{k+1}^2$  periods in expectation. Say that the  $k^{\text{th}}$  stage is *short* if it lasts for at most  $T(\varepsilon_{k+1})/\varepsilon_{k+1}^2$  periods, which is  $1/2k^2$  times the expected number. This event has probability at most  $1/k^2$ . Hence, given any positive integer  $k_0$ , the probability that a short stage occurs at some time after the  $k_0^{\text{th}}$  stage is at most  $\sum_{k>k_0} 1/k^2 \leq \int_{k_0}^{\infty} dx/x^2 = 1/k_0$ . If we let  $k_G^{**} = k_G^* \vee 16/\varepsilon$ , it follows that *the probability is at most  $\varepsilon/16$  that a short stage ever occurs after stage  $k_G^{**}$ .*

Now there exists a time  $T^{**}$  such that

$$\forall t \geq T^{**}, \quad P(N_t \geq k_G^{**} + 2) \geq 1 - \varepsilon/16. \quad (54)$$

We shall show that (53) holds for this value of  $T^{**}$ .

For each time  $t \geq T^{**}$ , define the event  $A_t$  to be the set of all realizations such that there is *at most one stage change between  $t - T(\varepsilon_{N_t})/\varepsilon_{N_t}^2$  and  $t$* , that is,

$$N_t \leq 1 + N_{t-T(\varepsilon_{N_t})/\varepsilon_{N_t}^2}. \quad (55)$$

Let  $A_t^c$  denote the complement of  $A_t$ . Since  $t \geq T^{**}$ , the probability is at least  $1 - \varepsilon/16$  that the process is at stage  $k_G^{**} + 2$  or higher at time  $t$ . Denote this event by  $B_t$ . If  $B_t$  and  $A_t^c$  both hold, then there were at least two stage changes between  $t - T(\varepsilon_{N_t})/\varepsilon_{N_t}^2$  and  $t$ , hence the *previous* stage change (before the current stage) was short. But we already know that the probability of a short stage at any time beyond stage  $k_G^{**}$  is at most  $\varepsilon/16$ . Hence  $P(A_t^c | B_t) \leq \varepsilon/16$  and  $P(B_t^c) \leq \varepsilon/16$ . Therefore

$$\forall t \geq T^{**}, \quad P(A_t^c) \leq P(A_t^c | B_t) + P(B_t^c) \leq 2(\varepsilon/16) = \varepsilon/8. \quad (56)$$

We now compute the probability that  $W_t < T(\varepsilon_{N_t})$ . By the preceding we know that

$$\begin{aligned} P(W_t < T(\varepsilon_{N_t})) &\leq P(W_t < T(\varepsilon_{N_t}) \mid A_t) + P(A_t^c) \\ &\leq P(W_t < T(\varepsilon_{N_t}) \mid A_t) + \varepsilon/8. \end{aligned} \quad (57)$$

Hence to establish (53) it suffices to show that

$$P(W_t < T(\varepsilon_{N_t}) \mid A_t) \leq 3\varepsilon/8. \quad (58)$$

Clearly,

$$\begin{aligned} P(W_t < T(\varepsilon_{N_t}) \mid A_t) &= \sum_k P(W_t < T(\varepsilon_{N_t}) \mid N_t = k, A_t) P(N_t = k \mid A_t) \\ &\leq \max_k P(W_t < T(\varepsilon_{N_t}) \mid N_t = k, A_t). \end{aligned} \quad (59)$$

Let  $N_t = k$ . The event  $A_t$  is the disjoint union of the event  $A_t^0$  in which no stage change occurs between  $t - T(\varepsilon_k)/\varepsilon_k^2$  and  $t$ , and the event  $A_t^1$  in which exactly one stage change occurs.

When  $A_t^0$  occurs,  $W_t \geq T(\varepsilon_k)/\varepsilon_k^2 > T(\varepsilon_k)$ , hence  $P(W_t < T(\varepsilon_k) \mid A_t^0) = 0$ . It remains only to show that  $P(W_t < T(\varepsilon_k) \mid A_t^1) \leq 3\varepsilon/8$ .

The conditional distribution of  $W_t$  is

$$f(w) \equiv P(W_t = w \mid N_t = k, A_t^1) = c_k (1 - p_k)^{T(\varepsilon_k)/\varepsilon_k^2 - w} p_k (1 - p_{k+1})^{w-1}, \quad (60)$$

where  $c_k$  is a positive constant and  $1 \leq w \leq T(\varepsilon_k)/\varepsilon_k^2$ . This follows because under  $A_t^1$  a single stage change occurs during the interval, and it occurs exactly  $W_t = w$  periods before period  $t$ . We may rewrite (60) in the form

$$f(w) = c'_k \left( \frac{1 - p_{k+1}}{1 - p_k} \right)^w \quad (61)$$

for some  $c'_k > 0$ . Since  $p_k > p_{k+1}$ ,  $f(w) \leq f(w + 1)$ . Hence for every  $T$  and  $w$  in the interval  $1 \leq T, w \leq T(\varepsilon_k)/\varepsilon_k^2$ ,

$$\sum_{w < T} f(w) \leq T f(T) \text{ and } \sum_{w \geq T} f(w) \geq (T(\varepsilon_k)/\varepsilon_k^2 - T) f(T). \quad (62)$$

In particular for  $T = T(\varepsilon_k)$  we have

$$\begin{aligned}
 P(W_t < T(\varepsilon_k)) &= \sum_{w < T(\varepsilon_k)} f(w) \\
 &= \frac{1}{1 + \frac{\sum_{w \geq T(\varepsilon_k)} f(w)}{\sum_{w < T(\varepsilon_k)} f(w)}} \leq \frac{1}{1 + \frac{T(\varepsilon_k)/\varepsilon_k^2 - T(\varepsilon_k)}{T(\varepsilon_k)}} = \varepsilon_k^2.
 \end{aligned} \tag{63}$$

Since  $t \geq T^{**}$ ,  $\varepsilon_{N_t} = \varepsilon_k \leq \varepsilon/4$  with probability at least  $1 - \varepsilon/4$ . Hence

$$P(W_t < T(\varepsilon_k)) \leq (\varepsilon/4)^2(1 - \varepsilon/4) + \varepsilon/4 < 3\varepsilon/8. \tag{64}$$

This establishes (53) and completes the proof of theorem 2.

## Appendix

Here we prove lemma 1, which is restated for easy reference.

**Lemma 1** *Let  $m = m_1 \vee m_2, \tau = \tau_1 \wedge \tau_2$  and  $\lambda = \lambda_1 \wedge \lambda_2$ , and suppose that  $0 < \lambda_i \leq \tau/8 \leq 1/8$  for  $i = 1, 2$ . There exist positive constants  $a, b$ , and  $c$  such that, for all  $t$ ,*

- i) *If state  $z_t = (p_t, q_t)$  is a  $(\tau_1/2, \tau_2/2)$ -equilibrium, a revision occurs at the end of period  $t$  with probability at most  $ae^{-bs}$  for all  $s$ .*
- ii) *If  $z_t$  is not a  $(2\tau_1, 2\tau_2)$ -equilibrium, each player revises at the end of period  $t$  with probability greater than  $1/2$  and both revise with probability greater than  $1/4$ , provided that  $s \geq c$ .*

*It suffices that  $a = 12m, b = \lambda\tau^2/256m$ , and  $c = 10^3m^2/\lambda\tau^2$ .*

**Proof:** The player's strategy revisions are triggered by the size of their realized regrets  $\widehat{r}_t^i$ . Hence we need to estimate the distribution of  $\widehat{r}_t^i$  conditional on the state at time  $t$ , namely,  $z_t = (p_t, q_t)$ . Recalling the definitions of  $\widehat{\alpha}_{j,t}^i$  and  $\widehat{\alpha}_t^i$  from Step 3 of regret testing, let

$$\alpha_{j,t}^i \equiv E(\widehat{\alpha}_{j,t}^i | (p_t, q_t)) \tag{A1}$$

and

$$\alpha_t^i \equiv E(\widehat{\alpha}_t^i | (p_t, q_t)). \tag{A2}$$

Recall that player 2 draws from his hat with probability  $1 - \lambda_2$ , and plays an action uniformly at random with probability  $\lambda_2$ . (The uniform distribution over actions when experimenting contrasts with the possibly non-uniform distribution over hats when a rejection occurs). Hence when player 1 chooses action  $j$  at time  $t$ , his expected payoff is

$$\alpha_{j,t}^1 = \sum_k ((1 - \lambda_2)(q_t)_k + \lambda_2/m_2)u_{j,k}^1.$$



Similarly, 1's expected payoff at time  $t$  is

$$\alpha_t^1 = \sum_{j,k} (p_t)_j ((1 - \lambda_2)(q_t)_k + \lambda_2/m_2) u_{j,k}^1.$$

Similar expressions hold for  $\alpha_{j,t}^2$  and  $\alpha_t^2$ . Define

$$r_t^i \equiv \max_j \alpha_{j,t}^i - \alpha_t^i. \quad (\text{A3})$$

Since  $E(\hat{\alpha}_{j,t}^i | (p_t, q_t)) = \alpha_{j,t}^i$  and  $E(\hat{\alpha}_t^i | (p_t, q_t)) = \alpha_t^i$  we can think of the difference,

$$\hat{r}_t^i = \max_j \hat{\alpha}_{j,t}^i - \hat{\alpha}_t^i,$$

as being an estimator of  $r_t^i$ .

Define the *estimation error in state*  $(p_t, q_t)$  to be

$$|\hat{r}_t^i - r_t^i|. \quad (\text{A4})$$

Next we estimate the distribution of the realized regret estimates  $\hat{r}_t^i$ .

**Claim:** If  $\lambda_i \leq 1/3$ , then for all  $\delta \leq 1/\sqrt{2m_i}$ , and for all times  $t$ ,

$$P(|\hat{r}_t^i - r_t^i| > \delta) \leq 6m_i e^{-\frac{s\lambda_i\delta^2}{16m_i}}. \quad (\text{A5})$$

**Proof:** Fix a player  $i$  and let  $(p_t, q_t)$  be the state on day  $t$ . Let  $N_{j,t}^i$  be the number of times action  $j$  is played on day  $t$  while player  $i$  is on the telephone. The average payoff during these times,  $\hat{\alpha}_{j,t}^i$ , is an average of  $N_{j,t}^i$  items, each of which is bounded between zero and one. By Azuma's inequality [1],

$$P(|\hat{\alpha}_{j,t}^i - \alpha_{j,t}^i| > \delta \mid (p_t, q_t), N_{j,t}^i) \leq 2e^{-N_{j,t}^i\delta^2/2}. \quad (\text{A6})$$

Let  $N_t^i = \sum_j N_{j,t}^i$ . The number of times  $i$  was not on the phone on day  $t$  is  $s - N_t^i$ , hence again by Azuma's inequality

$$P(|\hat{\alpha}_t^i - \alpha_t^i| > \delta \mid (p_t, q_t), N_t^i) \leq 2e^{-(s-N_t^i)\delta^2/2}. \quad (\text{A7})$$

Since for any two events  $\mathcal{A}$  and  $\mathcal{B}$ ,  $P(\mathcal{A} \cup \mathcal{B}) \leq P(\mathcal{A}) + P(\mathcal{B})$ , it follows from (A6) and (A7) that

$$P(|\widehat{r}_t^i - r_t^i| > 2\delta \mid (p_t, q_t), N_{1,t}^i, N_{2,t}^i, \dots, N_{m_i,t}^i) \leq 2 \sum_{j=1}^{m_i} e^{-N_{j,t}^i \delta^2/2} + 2e^{-(s-N_i^i)\delta^2/2}. \quad (\text{A8})$$

The next step is to estimate the size of the tail of the random variable  $N_{j,t}$ , which is binomially distributed  $B(\lambda_i/m_i, s)$ . We claim that:

$$P\left(\left|N_{j,t}^i - \frac{s\lambda_i}{m_i}\right| \geq \frac{s\lambda_i}{2m_i}\right) \leq 2e^{-s\lambda_i/20m_i}. \quad (\text{A10})$$

This can be derived from Bennett's inequality [3]. Consider a collection of  $n$  independent random variables  $U_1, \dots, U_n$  with  $\sup |U_i| \leq M$ ,  $EU_i = 0$ , and  $\sum_i EU_i^2 = 1$ . Then for every  $\tau > 0$ ,

$$P\left(\sum_i U_i \geq \tau\right) \leq \exp\left(\frac{\tau}{M} - \left(\frac{\tau}{M} + \frac{1}{M^2}\right) \log(1 + M\tau)\right). \quad (\text{A11})$$

We will apply this to the case of  $n$  i.i.d. random variables  $X_1, \dots, X_n$ , with  $\text{Var}(X_i) = \sigma^2$  and  $|X_i| \leq 1$ . Let  $U_i = (1/\sigma\sqrt{n})(X_i - EX)$ . Then  $|U_i| < 1/\sigma\sqrt{n}$ ,  $EU_i = 0$ , and  $\sum_{i=1}^n EU_i^2 = 1$ . Letting  $\tau = (\gamma/\sigma)\sqrt{n}$ ,  $M = 1/\sigma\sqrt{n}$ , and  $\bar{X} = \sum X_i/n$ , it follows from (A11) that

$$P(\bar{X} - EX \geq \gamma) \leq \exp(n\gamma - n(\gamma + \sigma^2) \log(1 + \gamma/\sigma^2)). \quad (\text{A12})$$

If we take  $\gamma = \sigma^2/2$  and use the fact that  $\log(3/2) \geq .4$ ,

$$P(\bar{X} - EX \geq \sigma^2/2) \leq \exp(-n\sigma^2/10). \quad (\text{A13})$$

When the  $X_i$ 's are binomial  $(p, n)$  with  $0 < p < .5$ , this implies

$$P(\bar{X} - p \geq p/2) \leq e^{-np/20}$$

and hence

$$P(|\bar{X} - p| \geq p/2) \leq 2e^{-np/20}, \quad (\text{A14})$$

from which (A10) follows immediately.

Consider the event  $\mathcal{B}$  in which all of the  $N_{j,t}^i$  lie within their expected value  $\lambda_i s/m_i$  plus or minus half their expected value:

$$\mathcal{B} \equiv \cap_j \{|N_{j,t}^i - \lambda_i s/m_i| \leq \lambda_i s/2m_i\}.$$

From (A10) it follows that the probability of the complementary event  $\mathcal{B}^c$  satisfies:

$$P(\mathcal{B}^c) \leq 2m_i e^{-s\lambda_i/20m_i} \quad (\text{A15})$$

Let  $\mathcal{A}$  be the event  $|\widehat{r}_t^i - r_t^i| > 2\delta$ . From (A8) we have

$$P(\mathcal{A}|\mathcal{B}) \leq 2m_i e^{-\lambda_i s\delta^2/4} + 2e^{-(s-N_t^i)\delta^2/2}. \quad (\text{A16})$$

But if  $\mathcal{B}$  holds, then  $s - N_t^i \geq s - 3\lambda_i s/2 = (1 - 3\lambda_i/2)s$ . By hypothesis  $\lambda_i \leq 1/3$ , hence  $s - N_t^i > s/2$  and

$$P(\mathcal{A}|\mathcal{B}) \leq 2m_i e^{-\lambda_i s\delta^2/4} + 2e^{-s\delta^2/4}. \quad (\text{A17})$$

Since  $P(\mathcal{A}) \leq P(\mathcal{A}|\mathcal{B}) + P(\mathcal{B}^c)$ , it follows from (A15) and (A17) that

$$\begin{aligned} P(\mathcal{A}) &= P(|\widehat{r}_t^i - r_t^i| > 2\delta) \leq 2m_i e^{-s\lambda_i\delta^2/4} + 2e^{-s\delta^2/4} + 2m_i e^{-s\lambda_i/20m_i} \\ &\leq 2(m_i + 1)e^{-s\lambda_i\delta^2/4} + 2m_i e^{-s\lambda_i/20m_i}. \end{aligned} \quad (\text{65})$$

Changing from  $\delta$  to  $\delta/2$  we obtain

$$P(|\widehat{r}_t^i - r_t^i| > \delta) \leq 2(m_i + 1)e^{\frac{-s\lambda_i\delta^2}{16}} + 2m_i e^{-\frac{s\lambda_i}{20m_i}}. \quad (\text{A19})$$

By assumption,  $\delta \leq 1\sqrt{2m_i}$ ; so  $1/20m_i \geq \delta^2/16 \geq \delta^2/16m_i$ . It follows that  $e^{-s\lambda_i\delta^2/16m_i} \geq e^{-s\lambda_i\delta^2/16} \geq e^{-s\lambda_i/20m_i}$ , hence (A19) implies

$$P(|\widehat{r}_t^i - r_t^i| > \delta) \leq (4m_i + 2)e^{\frac{-s\lambda_i\delta^2}{16m_i}} \leq 6m_i e^{\frac{-s\lambda_i\delta^2}{16m_i}}. \quad (\text{A20})$$

This establishes (A5) as claimed.

Recall that in state  $(p, q)$  the actual behavioral probabilities are, for player 1,

$$\widetilde{p} = (1 - \lambda_1)p + (\lambda_1/m_1)\vec{1}_{m_1},$$

and for player 2,

$$\tilde{q} = (1 - \lambda_2)q + (\lambda_2/m_2)\vec{1}_{m_2}. \quad (\text{A21})$$

If  $(p_t, q_t)$  is an  $(\varepsilon_1, \varepsilon_2)$ -equilibrium, the expected regrets  $r_t^i$  in state  $(p_t, q_t)$  satisfy the bound

$$r_t^i \leq \varepsilon_i + 2(\lambda_1 \vee \lambda_2). \quad (\text{A22})$$

(This follows from (A21) and the assumption that the payoffs are bounded between 0 and 1.)

To prove Lemma 1, part (i), assume that  $z_t = (p_t, q_t)$  is a  $(\tau_1/2, \tau_2/2)$ -equilibrium. Since by assumption,  $\lambda_1, \lambda_2 \leq \tau/8$ , where  $\tau = \tau_1 \wedge \tau_2$ , it follows from (A22) that

$$r_t^i \leq \tau_i/2 + \tau_i/4 = 3\tau_i/4. \quad (\text{A23})$$

In order for a rejection to occur, we must have  $\hat{r}_t^i > \tau_i$ , which by the preceding implies that  $|\hat{r}_t^i - r_t^i| > \tau_i/4$ . Letting  $\delta = \tau_i/4$ , it follows from (A5) that the probability of this occurring is less than  $6m_i e^{-\frac{s\lambda_i\tau_i^2}{256m_i}}$ . Thus the probability that one or both players reject is less than

$$\sum_{i=1}^2 6m_i e^{-\frac{s\lambda_i\tau_i^2}{256m_i}} \leq 12m e^{-s\lambda\tau^2/256m}.$$

This establishes lemma 1, part (i).

To prove part (ii) of the lemma, suppose that in state  $z_t$  at least one of the players, say  $i$ , can improve his payoff by more than  $2\tau_i$ . This implies  $r_t^i > 2\tau_i - \tau_i/4 \geq 7\tau_i/4$ . He rejects unless  $\hat{r}_t^i \leq \tau_i$ , which implies  $|\hat{r}_t^i - r_t^i| > 3\tau_i/4$ . By (A5) we know that the probability of this is less than  $6m_i e^{-\frac{9s\lambda_i\tau_i^2}{256m_i}} \leq 6m_i e^{-\frac{s\lambda_i\tau_i^2}{30m_i}}$ . Choose  $s$  large enough that

$$6m_i e^{-\frac{s\lambda_i\tau_i^2}{30m_i}} < 1/3.$$

This holds if

$$s > \frac{30m_i}{\lambda_i\tau_i^2} \ln(18m_i).$$

Noting that  $\ln x \leq x$  for  $x \geq 1$ , this simplifies to

$$s > \frac{540m_i^2}{\lambda_i\tau_i^2}.$$

Recalling that  $m = m_1 \vee m_2$ ,  $\tau_1 \wedge \tau_2$ , and  $\lambda = \lambda_1 \wedge \lambda_2$ , we see that this holds if  $s \geq c = 10^3 m^2 / \lambda \tau^2$ , as posited in the lemma. We have therefore shown that, if player  $i$  is out of equilibrium by more than  $2\tau_i$ , then  $i$  accepts with probability at most  $1/3$ , and rejects with probability at least  $2/3$ , which is certainly greater than  $1/2$ . If *both* players are in this situation (as posited in part (ii) of the lemma), then each revises with probability at least  $2/3$ . Since the union of these two events has probability at most 1, their intersection has probability at least  $1/3$  which is certainly greater than  $1/4$ . This concludes the proof of lemma 1 part ii).  $\square$

## 7 References

1. Azuma, K., Weighted sums of certain dependent random variables, *Tohoku Math. J.*, **19**, (1967), 357 - 367.
2. Bendor, Jonathan, Dilip Mookherjee, and Debraj Ray, Aspiration-based reinforcement learning in repeated interaction games: an overview, *International Journal of Game Theory*, **3**, (2001), 159-174.
3. Bennett, G. Probability inequalities for the sum of independent random variables, *JASA*, **57**, (1962), 33-45.
4. Börgers, Tilman, and Rajiv Sarin, Naïve reinforcement learning with endogenous aspirations, *International Economic Review*, **41**, (2000), 921-950.
5. Bush, R. R. and F. Mosteller, *Stochastic Models for Learning*. New York: John Wiley, 1955.
6. Cahn, Amotz, General procedures leading to correlated equilibria, *International Journal of Game Theory*, **33**, (2004), 21-40.
7. Cho, In-Koo, and Akihiko Matsui, Learning aspiration in repeated games, *Journal of Economic Theory*, forthcoming.
8. Erev, Ido, and Alvin E. Roth, Predicting how people play games: reinforcement learning in experimental games with unique, mixed strategy equilibria, *American Economic Review*, **88**, (1998), 848-881.
9. Foster, Dean P. and Rakesh Vohra, A randomization rule for selecting forecasts, *Operations Research*, **41**, (1993), 704-709.

10. Foster, Dean P., and Rakesh Vohra, Regret in the on-line decision problem, *Games and Economic Behavior*, **29**, (1999), 7-35.
11. Foster, Dean P., and H. Peyton Young, On the impossibility of predicting the behavior of rational agents, *Proceedings of the National Academy of Sciences of the USA*, **98**, no.222, (2001), 12848-12853.
12. Foster, Dean P., and H. Peyton Young, Learning, hypothesis testing, and Nash equilibrium, *Games and Economic Behavior*, **45**, (2003), 73-96.
13. Fudenberg, Drew, and David Levine, Consistency and cautious fictitious play, *Journal of Economic Dynamics and Control*, **19**, (1995), 1065-90.
14. Fudenberg, Drew, and David Levine, *The Theory of Learning in Games*, Cambridge MA: MIT Press, 1998.
15. Germano, Fabrizio, and Gabor Lugosi, Global convergence of Foster and Young's regret testing, Working paper, Departament d'Economia i Empresa, Universitat Pompeu Fabra, Barcelona, (2004).
16. Hart, Sergiu, and Andreu Mas-Colell, A simple adaptive procedure leading to correlated equilibrium, *Econometrica*, **68**, (2000), 1127-1150.
17. Hart, Sergiu, and Andreu Mas-Colell, A general class of adaptive strategies, *Journal of Economic Theory*, **98**, (2001), 26-54.
18. Hart, Sergiu, and Andreu Mas-Colell, Uncoupled dynamics do not lead to Nash equilibrium, *American Economic Review*, **93**, (2003), 1830-1836.
19. Hart, Sergiu, and Andreu Mas-Colell, Stochastic uncoupled dynamics and Nash equilibrium. Technical Report, Hebrew University of Jerusalem, 2004.

20. Jordan, James S., Bayesian learning in normal form games, *Games and Economic Behavior*, **5**, (1991), 368-386.
21. Jordan, James S., Three problems in learning mixed-strategy equilibria, *Games and Economic Behavior*, **5**, (1993), 368-386.
22. Kalai, Ehud, and Ehud Lehrer, Rational learning leads to Nash equilibrium, *Econometrica*, **61**, (1993), 1019-1045.
23. Karandikar, Rajeeva, Dilip Mookherjee, Debraj Ray, and Fernando Vega-Redondo, Evolving aspirations and cooperation, *Journal of Economic Theory*, **80**, (1998), 292-331.
24. Karlin, Samuel, and H.M. Taylor, A First Course in Stochastic Processes, 2nd edition. New York: Academic Press, 1975.
25. Nachbar, John H., Prediction, optimization, and learning in games, *Econometrica*, **65**, (1997), 275-309.
26. Nachbar, John H., Beliefs in repeated games, *Econometrica*, **73**, (2005), 459-480.
27. Young, H. Peyton, Strategic Learning and Its Limits. Arne Ryde Memorial Lectures. Oxford, UK: Oxford University Press, 2004.