

Learning by Forgetful Players

SJAAK HURKENS*

CentER, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, The Netherlands

Received February 28, 1994

A product set of pure strategies is said to be closed under best replies if all best replies against all possible mixtures of these strategies are contained in the set. Minimal sets with this property are called minimal curb sets. This paper presents a dynamic learning process that has two main characteristics: Players have a bounded memory and they play best replies against beliefs, formed on the basis of strategies used in the recent past. It is shown that this learning process leads the players to playing strategies from a minimal curb set. Several variations of the process are considered. *Journal of Economic Literature* Classification Numbers: C70, C72. © 1995 Academic Press, Inc.

1. INTRODUCTION

A product set of pure strategies is said to be closed under best replies if all best replies against all possible mixtures of these strategies are contained in the set. Minimal sets with this property are called minimal curb sets (Basu and Weibull, 1991). As we will see later, curb sets are closely related to the better known persistent retracts. Kalai and Samet (1984) showed that every game has at least one persistent retract and that every persistent retract contains at least one (proper) Nash equilibrium. This enabled them to introduce the persistent equilibrium as a refinement of the Nash equilibrium concept.

Both concepts have been used in the literature. Kalai and Samet (1985) used persistency to achieve efficiency in unanimity games that are repeated

* I thank Peter Borm, Eric van Damme, Jürgen Eichberger, Drew Fudenberg, participants of the IIASA workshop on Evolutionary Game Dynamics in Biology and Economics, and two anonymous referees for helpful discussions and comments. This research was sponsored by the Foundation for the Promotion of Research in Economic Sciences, which is part of the Netherlands Organization for Scientific Research (NWO).

as long as no agreement is reached. Blume (1993a) used the persistent retract as a set-valued solution concept in sender receiver games. Blume (1993b) shows that, in one-sided cheap talk games, equilibria in minimal curb sets sometimes select the sender's preferred outcome. Hurkens (1993) shows that, in games where several players have the possibility to send costly messages, minimal curb sets always select the outcome preferred by all senders. Van Damme and Hurkens (1993) applied the concepts of curb and persistency in games of endogenous timing and Balkenborg (1993) did so in finitely repeated games.

In most of these papers it is argued informally that the concepts of curb and persistency have a dynamic and evolutionary flavor. However, few or no attempts have been made to support this idea with an evolutionary foundation of the concepts.

We construct a dynamic learning process to support these concepts. Roughly speaking, the learning evolves as follows: A particular game is played at discrete points in time. For each role in this game there is a pool of players. At the beginning of each period one player is drawn from each pool. These players will play the game in that period. Players have a bounded memory. On the basis of strategies played in the recent past, they form expectations about the strategies the other players will use and best respond to these expectations. We assume that different players within the same pool may have different beliefs and therefore they may choose different actions. It is shown that, if the memory is long enough, play will settle down in a minimal curb set.

In some respects our results are stronger than those obtained thus far in the literature on learning. First, in contrast to Young (1993) we do not need to restrict attention to a special class of games. Second, the set of curb strategies is a subset of the set of rationalizable strategies (Bernheim, 1984; Pearce, 1984). Hence, our learning process reduces the number of "plausible" strategies. This is in contrast with Milgrom and Roberts (1991), who show that a sequence that is consistent with adaptive learning will eventually lie within the set of serially undominated strategies, which is a superset of the set of rationalizable strategies. In the final section we show that it is the forgetfulness of the players that accounts for this difference.

From the main and basic theorem we derive several results for learning processes where players learn in a somewhat different way. Play still settles down in minimal curb sets when some players do not play best responses to past play, but are more sophisticated than that or, on the contrary, are less sophisticated. If we allow players to have beliefs as if the other players in the game correlate their actions, play settles down in a primitive formation (Harsanyi and Selten, 1988), a variant of a minimal curb set. When players are uncertain, the process does not converge to a curb set but to related solution concepts as curb*, robust, or persistent sets, depending on how

the uncertainty is modeled. The learning processes presented in this paper may give the reader some insight in the differences and similarities between these related concepts. We also characterize two classes of games where our results go through, even if the players only observe the outcomes of past play, instead of the strategies.

The rest of the paper is organized as follows. In Section 2 we introduce some preliminaries concerning Markov chains and curb sets. Section 3 describes the model of learning as a Markov chain. Section 4 contains the main result: the ergodic sets of the Markov chain correspond one-to-one to the minimal curb sets of the underlying game. In Section 5 and 6 the above-mentioned variations of the learning process are considered. In Section 7 we consider the possibility that players make mistakes with small probability. Section 8 compares the present paper to Milgrom and Roberts (1991).

2. PRELIMINARIES

Let $G = (S_1, \dots, S_n, u_1, \dots, u_n)$ be a finite game with player set $N = \{1, \dots, n\}$. Let $S = \prod_{i=1}^n S_i$ and $S_{-i} = \prod_{j \neq i} S_j$. For any finite set X let $\Delta(X)$ denote the set of probability distributions over X . For a distribution $\mu \in \Delta(S)$ let $\mu_i \in \Delta(S_i)$ be the marginal on S_i , and let $\mu_{-i} \in \Delta(S_{-i})$ be the marginal on S_{-i} , i.e.,

$$\begin{aligned}\mu_i(s_i) &= \sum_{s_{-i} \in S_{-i}} \mu(s_i, s_{-i}) & (s_i \in S_i) \\ \mu_{-i}(s_{-i}) &= \sum_{s_i \in S_i} \mu(s_i, s_{-i}) & (s_{-i} \in S_{-i}).\end{aligned}$$

Of special interest are the probability distributions whose marginals on S_1, \dots, S_n are independent. The sets of these probability distributions will be denoted by Σ and Σ_{-i} , respectively. Although they are formally not the same, we will identify Σ with $\prod_{i=1}^n \Delta(S_i)$ and Σ_{-i} with $\prod_{j \neq i} \Delta(S_j)$ and trust that no confusion will result.

For $\mu \in \Delta(S)$ and $i \in N$ we let $\text{BR}_i(\mu_{-i})$ denote the set of pure best replies against μ_{-i} . Let $\text{BR}(\mu) = \prod_{i=1}^n \text{BR}_i(\mu_{-i})$. For $F \subset \Delta(S)$ let $\text{BR}_i(F) = \bigcup_{\mu \in F} \text{BR}_i(\mu_{-i})$ and $\text{BR}(F) = \prod_{i=1}^n \text{BR}_i(F)$.

DEFINITION 1. A non-empty cartesian product set $C = \prod_{i=1}^n C_i \subset S$ is said to be closed under best replies (or C is a curb set) if $\text{BR}(\prod_{i=1}^n \Delta(C_i)) \subset C$. Such a set is called a minimal curb set if it does not properly contain a curb set. Strategies contained in minimal curb sets are called curb strategies.

It is straightforward to show that $BR(\prod_{i=1}^n \Delta(C_i)) = C$ for any minimal curb set C . The notion of curb sets was introduced by Basu and Weibull (1991). Curb is mnemonic for closed under rational behavior.

A strict Nash equilibrium is a curb set as a singleton. Strict Nash equilibria have almost all desired properties one can hope for, except existence. Many of these properties carry over to minimal curb sets. For instance, every curb set contains the support of a proper equilibrium (Kalai and Samet, 1984; Balkenborg, 1992). Moreover, every game has at least one minimal curb set since S is curb.

Minimal curb sets can be viewed as a set-valued generalization of strict equilibria: when an outsider recommends to all players to play strategies from a minimal curb set C , then all players will follow this recommendation if they expect the other players to do so. The comparison with strict equilibria is not completely valid: minimal curb sets may contain weakly dominated strategies.

Before we go further let us consider some examples where minimal curb sets have some cutting power.

EXAMPLE A. Let G be given by the normal form in Fig. 1a. This is a pure coordination game. Since (T, L) and (B, R) are strict equilibria it is easy to see that $\{(T, L)\}$ and $\{(B, R)\}$ are minimal curb sets and that there are no other ones. In particular, the support of the mixed equilibrium is not contained in any minimal curb set.

EXAMPLE B. Now consider the game in which player 1 has the choice between playing the game from Fig. 1a and an outside option O , yielding both players a payoff of 3. The normal form representation of this game is given in Fig. 1b. This game has a unique minimal curb set, namely $\{(T, L)\}$.

These two examples are nice because the minimal curb sets are singletons and hence consist of one strict Nash equilibrium. In the following example, in contrast to those above, the unique minimal curb set is not a singleton.

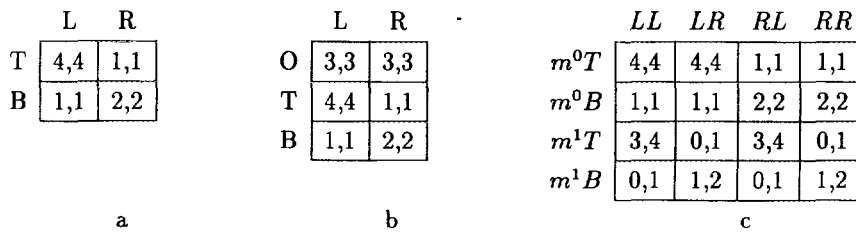


FIGURE 1

EXAMPLE C. Suppose that player 1 can send one of two messages, m^0 or m^1 , to player 2 before the game from Fig. 1a is played. Suppose that it costs player 1 i units to send m^i . Let ma denote player 1's strategy "I send message m and choose action a " and let a^0a^1 denote player 2's strategy "I choose action a^i if I receive message m^i ." Then the (reduced) normal form of the game with pre-play communication is given in Fig. 1c. Now it can be checked that $\{m^0T\} \times \{LL, LR\}$ is the unique minimal curb set of this extended game. The set is not a singleton but consists only of equilibria that involve sending the cheapest message and then playing the equilibrium preferred by player 1. In Hurkens (1993) similar results are obtained for a whole class of games with n players among which k have the possibility to send a costly message.

In the next section we will describe the learning process by means of a Markov chain. Therefore we will need some basic notions from the theory of Markov chains.

A finite stationary Markov chain is characterized by a pair (X, P) , where X is a finite state space and $P: X \times X \rightarrow [0, 1]$ is a transition matrix. The interpretation is that $P(x, x')$ is the probability that the process will move from x to x' in one period. We will denote $x \rightsquigarrow x'$ if there exist $k \in \mathbb{N} \cup \{0\}$, $x_0, \dots, x_k \in X$ with $x_0 = x$, $x_k = x'$ and $P(x_i, x_{i+1}) > 0$ ($i = 0, \dots, k - 1$). Now \rightsquigarrow defines a weak order on X . Hence, we can define an equivalence relation on X :

$$x \sim y \Leftrightarrow x \rightsquigarrow y \text{ and } y \rightsquigarrow x.$$

Let $[x]$ denote the equivalence class that contains x and let $Q = \{[x] | x \in X\}$ denote the set of equivalence classes. We define a partial order \leq on Q :

$$[x] \leq [y] \Leftrightarrow y \rightsquigarrow x.$$

The minimal elements with respect to the order \leq are called *ergodic sets*. The other elements are called *transient sets*. If the process leaves a transient set it can never return to that set. And if the process is in an ergodic set it can never leave this set. The elements of these sets are called *ergodic* and *transient states*, respectively. We have the following theorem.

THEOREM 1. *In any finite Markov chain, no matter where the process starts, the probability after k steps that the process is in an ergodic state tends to 1 as k tends to infinity.*

Proof. See, e.g., Kemeny and Snell (1976).

3. THE LEARNING PROCESS

According to the Bayesian approach, a player forms some expectation about the strategies that will be played by the other players and best responds to his expectation. How these expectations are formed is not clear. When the same game has been played before, possibly by different people, it seems reasonable to suggest that expectations are formed on the basis of information on past play. One way of using this information is to assume that a player's belief corresponds to the empirical frequency of strategies used in the past. This way of forming beliefs, known as fictitious play (Brown, 1951; Robinson, 1951), perhaps makes sense in matching models, but it is certainly not the only possible way of forming beliefs. One drawback of fictitious play is that it assumes that all people always form expectations in the same way. This implies that if different people have the same information, they will form the same beliefs and consequently they choose the same action. One can create some stochastic variability in the process by assuming that people only draw an incomplete sample of the information, as in Young (1993). There it is assumed that players learn how the game was played in m out of the most recent K periods. The players use a fictitious play rule to map samples into beliefs, and best respond to these beliefs. The great technical advantage of Young's approach is that the learning process can be described by a finite Markov chain on the state space $H = S^K$, consisting of all sequences of length K drawn from S . In order to determine the ergodic sets of such Markov chains, one needs only to specify which transitions occur with positive probability and which occur with zero probability.

We will also describe a learning process by means of a finite Markov chain, but we allow more variability in the responses of the players. In fact, we allow the degree of variability that is present in Milgrom and Roberts' (1991) definition of adaptive play.¹

Let $G = (S, u)$ be an n -person normal form game. Fix a positive integer K . Suppose we have a finite population of individuals that is partitioned into non-empty classes V_1, \dots, V_n . The members of V_i are candidates to play role i in the game, and they all have the same payoff function u_i . Let $t = 0, 1, 2, \dots$ denote successive time periods. Game G is played once every period. In period t one individual is drawn from each class V_i . These individuals are going to play the appropriate roles in the game this period. We will refer to the individual that is drawn from V_i to play the game in the current period as player i , although the identity of this player may change from time to time. Player i receives some, but not necessarily all, information about play in the recent K periods. Then he chooses a pure

¹ See Section 8 for a comparison between the present paper and Milgrom and Roberts (1991).

strategy according to some rule. We will define below what kind of information a player may receive, and how he chooses a strategy as a function of this information. Then the players are put back in their class. This ends period t and we move up to period $t + 1$.

Since we will assume that all the rules are time-independent, this learning process can be described by a stationary Markov chain on the state space $H = S^K$. Call $\hat{h} \in H$ a *successor* of $h \in H$ if \hat{h} is obtained from h by deleting the leftmost element and by adding some element $s \in S$ to the right. Let $r(\hat{h})$ denote the rightmost element of $\hat{h} \in H$. For $h = (s^{-K}, \dots, s^{-1}) \in H$ let $\pi_i(h) = \{s_i^{-K}, \dots, s_i^{-1}\}$ denote the set of strategies played by player i in the recent past. We will assume that our learning process is described by a transition matrix $P \in \mathcal{P}$, where \mathcal{P} is defined as follows.²

DEFINITION 2. Let \mathcal{P} denote the set of transition matrices P that satisfy for all histories $h, \hat{h} \in H$,

$$P(h, \hat{h}) > 0 \Leftrightarrow \begin{cases} \hat{h} \text{ is a successor of } h, \text{ and} \\ r_i(\hat{h}) \in \text{BR}_i(\mu^i) \text{ for some } \mu^i \in \prod_{j \neq i} \Delta(\pi_j(h)) \end{cases} \quad (\text{all } i).$$

We will give two interpretations of a learning process that is described by some $P \in \mathcal{P}$. The first interpretation is close to the model of Young (1993). Fix a positive integer L . Before player i chooses a strategy in period t , he receives information about how the game was played by player j in the recent past, for all $j \neq i$. He receives L draws *with* replacement from the set $\{s_j(t - K), \dots, s_j(t - 1)\}$. A way of thinking about this sampling procedure is that player i passively hears about L precedents concerning the way player j played the game before. But player i is unaware of the fact that he might hear about the same precedent several times. Assume that all draws are independent, but more importantly, assume that each combination of draws occurs with positive probability. Player i 's belief about the behavior of player j corresponds to the empirical frequency of strategies in the sample of size L . Hence, this belief is one of a finite number of possible probability distributions. Namely, let $h = (s(t - K), \dots, s(t - 1))$ denote the recent history and let $\pi_j(h) = \{s_j(t - K), \dots, s_j(t - 1)\}$ denote the set of strategies played by player j in the recent past. Now player i 's belief about player j 's behavior is contained in the set

$$B_j(h, L) = \{\mu_j \in \Delta(\pi_j(h)) \mid \mu_j(s_j) = l/L \text{ for some } l \in \{0, 1, \dots, L\}\}.$$

² A transition matrix describes a learning process for a fixed game, G , and a fixed length of the memory, K . We will, however, suppress superscripts G and K .

We call the set $B^i(h, L) = \prod_{j \neq i} B_j(h, L)$ the L -grid distribution space for i induced by h . Note that as L increases, the grid becomes finer and finer, and $B^i(h, L)$ “approaches” $\prod_{j \neq i} \Delta(\pi_j(h))$. There exists a “generic” class of games for which it suffices, for the purpose of this paper, to choose L sufficiently large. However, in general we need a little bit more and therefore we assume that our learning process is described by some $P \in \mathcal{P}$.

Another interpretation of a learning process that is described by a transition matrix $P \in \mathcal{P}$ is the following. Suppose that the individuals in a class have different personal characteristics: They use the information on past play to know which strategies will certainly not be used (namely the ones that have not been played in the recent history). But each individual makes his own personal assessment of the probabilities with which the remaining strategies will be played. Some people are very optimistic and expect the best, while others are very pessimistic and expect the worst. And there will be a lot who have some intermediate beliefs. Of course, we need sufficient diversity in the different classes when this learning process is to be described by some $P \in \mathcal{P}$. Note, however, that this does not necessarily mean that these classes are large. Suppose that for each strategy $s_i \in S_i$, there is some individual in V_i who plays s_i whenever it is a best reply to some belief that puts positive weight only on strategies that were played recently. (And he chooses a best reply to the most recent strategy otherwise.) Then we only need $|S_i|$ individuals in class V_i .

In the next section we will state and prove the main theorem of this paper: Play will settle down in a minimal curb set.

4. ERGODIC SETS

Fix $K \in \mathbb{N}$ as the length of the histories. Recall from Section 2 that $h \rightsquigarrow \hat{h}$ means that there exist $k \in \mathbb{N}$, $h^0, \dots, h^k \in H = S^K$ such that $h^0 = h$, $h^k = \hat{h}$, and $P(h^i, h^{i+1}) > 0$. Now \rightsquigarrow defines a weak order on H and hence we can define an equivalence relation on H and an order on the set of equivalence classes of H . We will be interested in the minimal elements of this order, the ergodic sets.

Let C be a minimal curb set of $G = (S, u)$. We say that $h \in H$ is a C -history if $h \in C^K$. We call h a *curb history* if it is a C -history for some minimal curb set C .

Now we are ready to state the main theorem.

THEOREM 2. *There exists $\underline{K} \in \mathbb{N}$ such that for all $K \geq \underline{K}$ and every Markov chain with a transition matrix $P \in \mathcal{P}$:*

- (i) If $Z \subset H$ is an ergodic set then $Z \subset C^K$ for some minimal curb set C .
- (ii) For every minimal curb set C there exists exactly one subset $Z \subset C^K$ that is ergodic.
- (iii) For each minimal curb set C and each strategy $\bar{s} \in C$ there exists an ergodic state h with $r(h) = \bar{s}$.

This theorem states that, if the history is long enough, any ergodic set is a set of C -histories for some minimal curb set C and that the set of C -histories contains one ergodic set. Hence, the ergodic states are curb histories. Moreover, once the ergodic set contained in C^K is entered, every strategy $\bar{s} \in C$ is played infinitely often. From Theorem 1 then the following corollary follows.

COROLLARY 1. *The probability that the players are playing a curb strategy profile after k steps of the learning process tends to 1 as k tends to infinity if histories are sufficiently long.*

The intuition for the theorem is quite clear. By having a large enough memory, players may have beliefs with large supports. This means that best replies against all kinds of mixtures will be played now and then. This creates so much stochastic variability that players sooner or later will play curb strategies. When they keep drawing the “right” samples, they will keep best responding against curb strategies and hence will play curb strategies again. It might happen that they will do this K periods in a row. The probability that this happens at a specific point in time is only small, but with probability one it will happen eventually. By that time all non-curb strategies will be forgotten. The strategies that will be played from that point on will depend on the sample drawn, but it is sure that it will be curb strategies again.

Before we give a formal proof we make two remarks about Theorem 2. First, note that assertions (i) and (ii) do not imply that C^K is an ergodic set whenever C is a minimal curb set. Still, the reader may think that the only ergodic set contained in C^K is C^K itself. However, the following example shows that C^K need not be ergodic.

Consider the game in Fig. 2. This game has only one curb set, namely

	a_2	b_2	c_2
a_1	4,1	1,4	2,3
b_1	1,4	4,1	2,3
c_1	3,2	3,2	0,0

FIGURE 2

the set of all pure strategy combinations. But the profile $\bar{h} = (c, c, \dots, c)$ cannot be reached under the learning process from any other history. This is so because c is only a best reply against some mixtures of a and b . Hence, there exists no h with $P(h, \bar{h}) > 0$ and \bar{h} is not contained in the ergodic set.

The second remark concerns the length of the histories. In the proof of Theorem 2 we will use a lower bound on K , but that bound is not tight. On the other hand, it is important that histories are not too short, as the example from Fig. 3 shows.

It is not difficult to see that if $K = 2$, then the set of histories $\{(s^{-2}, s^{-1}) | s^{-j} \in \{T, M, B\} \times \{l, c, r\}\}$ contains an ergodic set. Take for example the history (Tl, Mr) . Agents from pool V_1 will play a best reply against $\alpha l + (1 - \alpha)r$, for some $\alpha \in [0, 1]$. Hence, they will play T or B . But the unique minimal curb set is the singleton $\{(A, a)\}$. So the history must not be too short. Note that if $K = 3$ and the process is in state (Tl, Mr, Mc) , then there will be some agent in V_1 who will play A , since A is the best reply against $\frac{1}{3}l + \frac{1}{3}c + \frac{1}{3}r$.

Note that the game shown in Fig. 3 has a unique equilibrium, namely (A, a) . This equilibrium is strict. Since every curb set contains the support of a Nash equilibrium and since a strict equilibrium forms a curb set as a singleton, it follows that this game has a unique minimal curb set. Hence, if players behave as described by our learning process then they will eventually play the equilibrium. This reasoning holds for all games that have a unique equilibrium that happens to be strict. So we proved

COROLLARY 2. *Suppose that s is the unique Nash equilibrium of G and that s is strict. The probability that players are playing the equilibrium after k steps of the learning process tends to 1 as k tends to infinity if histories are sufficiently long.*

The remainder of this section contains a formal proof of Theorem 2. First we introduce some notation and state a lemma.

Let F be a non-empty subset of S . We define the projection of F on S_i

	a	l	c	r
A	4,4	2,2	2,2	2,2
T	2,2	5,0	0,5	0,0
M	2,2	0,0	5,0	0,5
B	2,2	0,5	0,0	5,0

FIGURE 3

as $p_i(F) = \{f_i | f \in F\}$ and we define $\text{span}(F) = \prod_{i=1}^n p_i(F)$. Hence, $\text{span}(F)$ is the smallest Cartesian product set in S that contains F . Similarly, for a history $h = (s^{-K}, \dots, s^{-1})$ we define $\pi_i(h) = \{s_i^{-K}, \dots, s_i^{-1}\}$ and $\text{span}(h) = \prod_{i=1}^n \pi_i(h)$. We say that $B \subset S$ spans F if $\text{span}(B) = \text{span}(F)$.

For a history h let $\mathcal{B}^{\text{ind}}(h) = \{\mu \in \Sigma | \text{supp}(\mu) \subset \text{span}(h)\}$. This set contains all independent beliefs a Bayesian player might have when the process is in state h . Similarly, we define for a set $F \subset S$, $\mathcal{B}^{\text{ind}}(F) = \{\mu \in \Sigma | \text{supp}(\mu) \subset \text{span}(F)\}$. Let $M = \max_i |S_i|$.

LEMMA 1. *Let $s^1, \dots, s^T \in S$ be such that $s^{t+1} \notin \text{span}(\{s^1, \dots, s^t\})$ for all $t = 1, \dots, T - 1$. Then $T \leq \sum_{i=1}^n |S_i| - (n - 1)$.*

Proof. Easy and hence omitted. ■

Proof of Theorem 2. Take $\underline{K} = \sum_{i=1}^n |S_i| - (n - 1) + M$ and let $K \geq \underline{K}$. Let $P \in \mathcal{P}$.

Let $h^t = (x^{K-t}, \dots, x^1, s^1, \dots, s^t)$ be a particular history and assume that $F^t = \text{span}(\{s^1, \dots, s^t\})$ is not a curb set. Then there exists $s^{t+1} \in \text{BR}(\mathcal{B}^{\text{ind}}(F^t)) \setminus F^t$. Let $h^{t+1} = (x^{K-t-1}, \dots, x^1, s^1, \dots, s^{t+1})$. Then $P(h^t, h^{t+1}) > 0$. Starting from an arbitrary history h^1 we can apply this argument repeatedly. By Lemma 1 we know that there exists $T \leq \underline{K} - M$ such that $h^1 \rightsquigarrow h^T = (x^{K-T}, \dots, x^1, s^1, \dots, s^T)$ and such that $F^T = \text{span}(\{s^1, \dots, s^T\})$ is a curb set. Let $C \subset F^T$ be a minimal curb set and let $\{b^1, \dots, b^M\}$ span C . Since every strategy in a minimal curb set is a best reply to some belief concentrated on this set and since $K \geq M + T$, we have $h^T \rightsquigarrow (\dots, s^1, \dots, s^T, b^1, \dots, b^M) \rightsquigarrow (b^1, \dots, b^M, b^M, \dots, b^M)$.

The above shows that for every history h , there exists a minimal curb set C such that for every set $\{b^1, \dots, b^M\}$ that spans C , we have $h \rightsquigarrow (b^1, \dots, b^M, b^M, \dots, b^M)$. Furthermore, the definition of \mathcal{P} implies that if h is a C -history and $h \rightsquigarrow \hat{h}$, then \hat{h} is also a C -history.

The second observation implies that the set of C -histories contains an ergodic set, for any minimal curb set C . The first observation then implies that the set of C -histories contains exactly one ergodic set and that there are no other ergodic sets. Assertion (iii) follows from the observation that the spanning set $\{b^1, \dots, b^M\}$ can be chosen such that $b^M = \bar{s}$. ■

5. VARIATIONS ON THE SAME THEME

We remarked before that one only needs to know which entries of the transition matrix are positive and which are zero in order to characterize the ergodic sets. In the proof of Theorem 2 we used that certain entries are positive (together with Lemma 1) to show that the process can move from any history h to a curb history \hat{h} in a finite number of periods.

Furthermore, we used the fact that certain entries are zero to ensure that the process cannot leave the set of C -histories for any curb set C .

It is possible to prove Theorem 2 for an even bigger class of transition matrices. Let $P \in \mathcal{P}$ and let \tilde{P} be a transition matrix that satisfies, for any minimal curb set C ,

$$P(h, \hat{h}) > 0 \Rightarrow \tilde{P}(h, \hat{h}) > 0 \tag{5.1}$$

$$h \in C^K \text{ and } \tilde{P}(h, \hat{h}) > 0 \Rightarrow \hat{h} \in C^K. \tag{5.2}$$

Let $\tilde{\mathcal{P}}$ denote the set of all such transition matrices. It is obvious that Theorem 2 holds for all $P \in \tilde{\mathcal{P}}$. We will consider two subsets of $\tilde{\mathcal{P}}$, namely $\mathcal{P}^{\text{soph}}$ and \mathcal{P}^{mim} . The transition matrices in these sets correspond to learning processes where some players are more sophisticated (in the case of $\mathcal{P}^{\text{soph}}$) or less sophisticated (in the case of \mathcal{P}^{mim}). It turns out that for these two classes we can prove slightly stronger results.

5.1. More and Less Sophisticated Players

Suppose that not all individuals in the classes are Bayesian players, but that some individuals are mimickers. Mimickers do not form expectations but just observe how other agents in the same role have played the game during (some of) the last K periods. Then they choose one of these strategies at random. When we retain our assumption about the Bayesian players, this learning process can be described by a transition matrix $P \in \mathcal{P}^{\text{mim}}$, where \mathcal{P}^{mim} is defined as follows.

DEFINITION 3. Let \mathcal{P}^{mim} denote the set of transition matrices P that satisfy for all histories $h, \hat{h} \in H$,

$$P(h, \hat{h}) > 0 \Leftrightarrow \begin{cases} \hat{h} \text{ is a successor of } h, \text{ and} \\ r_i(\hat{h}) \in \text{BR}_i(\mathcal{B}^{\text{ind}}(h)) \text{ or } r_i(\hat{h}) \in \pi_i(h) \end{cases} \quad (\text{all } i).$$

Obviously, $\mathcal{P}^{\text{mim}} \subset \tilde{\mathcal{P}}$, and hence Theorem 2 holds for all $P \in \mathcal{P}^{\text{mim}}$. We can prove a slightly stronger result: All curb histories are ergodic states.

THEOREM 3. *There exists $\underline{K} \in \mathbb{N}$ such that for all $K \geq \underline{K}$ and for every Markov chain with a transition matrix $P \in \mathcal{P}^{\text{mim}}$, $Z \subset H$ is an ergodic set if and only if $Z = C^K$ for some minimal curb set C .*

Proof. Using the proof of Theorem 2, it suffices to show that if C is a minimal curb set and h and \hat{h} are C -histories, then $h \rightsquigarrow \hat{h}$.

Let $\hat{h} = (s^{-K}, \dots, s^{-1})$. We can choose a set $B = \{b^1, \dots, b^M\}$ that spans C such that $s^{-j} \in \text{span}(\{b^{M-j+1}, \dots, b^M\})$, for $j = 1, \dots, M$. From the proof

of Theorem 2 we know that $h \rightsquigarrow (b^1, \dots, b^M, s^{-K}, \dots, s^{-(M+1)}) =: \bar{h}$. Because of the special way we chose B (and because players sometimes mimic) we have $\bar{h} \rightsquigarrow \hat{h}$. ■

It is possible to prove Theorem 3 with a smaller lower bound on the length of the memory by making full use of the presence of the mimickers. We will not pursue that here. We just remark that for weakly acyclic games, the class of games considered in Young (1993), we could take $K = 1$.

The learning process we considered implies that Bayesian players play best responses against past play. If a player knew that other players are following this process, he could do better by playing a strategy that is a best reply against a strategy profile, consisting of best responses for the other players against past play. Of course, we may have players who foresee that others are going to play best responses to best replies to past play. We could have even more sophisticated players. When we assume that in a class many different levels of sophistication are represented, we have a learning process with sophisticated players. (See also Milgrom and Roberts, 1991.)

Formally, let h be a particular history and let $T^0(h) = \text{span}(h)$. Define recursively $T^{j+1}(h) = \text{span}(T^j(h) \cup \text{BR}(\mathcal{B}^{\text{ind}}(T^j(h))))$. Since $T^{j+1}(h) \supset T^j(h)$ and S is finite, $T^\infty(h) = \text{span}(\bigcup_{j=0}^\infty T^j(h))$ is well defined. Again, we define a whole set of transition matrices that correspond to learning processes with sophisticated players. We will denote this class by $\mathcal{P}^{\text{soph}}$, where $\mathcal{P}^{\text{soph}}$ is defined as follows.

DEFINITION 4. Let $\mathcal{P}^{\text{soph}}$ denote the set of transition matrices P that satisfy for all histories $h, \hat{h} \in H$,

$$P(h, \hat{h}) > 0 \Leftrightarrow \begin{cases} \hat{h} \text{ is a successor of } h, \text{ and} \\ r(\hat{h}) \in \text{BR}(\mathcal{B}^{\text{ind}}(T^\infty(h))). \end{cases}$$

It is obvious that $\mathcal{P}^{\text{soph}} \subset \tilde{\mathcal{P}}$ and hence Theorem 2 is valid also for this class. We can prove a stronger result: In the presence of sophisticated players we need only a memory of length one. The intuition for this result is that sophisticated players can do all the learning in their heads. They might foresee all the steps that needed to be executed in the case of no sophisticated players.

THEOREM 4. For all $K \geq 1$ and all Markov chains with a transition matrix $P \in \mathcal{P}^{\text{soph}}$ we have $Z \subset H$ is an ergodic set if and only if $Z = C^K$ for some minimal curb set C .

Proof. For notational convenience we just give the proof for $K = 1$. Now $H = S$ and we can define $T^\infty(s)$ for all $s \in S$. Note that $T^\infty(s)$ is a

curb set and hence there exists a minimal curb set $\bar{C} \subset T^\infty(s)$. If $\bar{s} \in \bar{C}$ then $P(s, \bar{s}) > 0$.

Note that if $s \in C$ for some minimal curb set C then $T^\infty(s) = C$. Hence, if $s, \bar{s} \in C$, then $P(s, \bar{s}) > 0$. ■

The reader may have noticed that this sophisticated learning process has some similarities with the notion of rationalizability (Bernheim, 1984; Pearce, 1984). The difference is that rationalizability corresponds with a process of iterative elimination of strategies that are never best replies (starting with the whole space of strategy profiles) whereas our learning process implies the addition of best replies (starting from a history). The bounded memory of the players causes play to settle down in a *minimal* curb set.

The similarity of rationalizable and curb strategies has already been pointed out by Basu and Weibull (1991) and Balkenborg (1992): Call a set $C = \prod_{i=1}^n C_i$ *tight* if $\text{BR}(\prod_{i=1}^n \Delta(C_i)) = C$. The maximal tight set is the set of rationalizable strategies, while the minimal tight sets are just the minimal curb sets. In particular, every curb strategy is rationalizable.

5.2. Uncertain Players

Consider the game shown in Fig. 4. This game has a unique curb set: it consists of all pure strategy profiles. When players behave as described by any of the learning processes they will regularly be playing (B, R) ! This might seem a bit strange. It could not happen if the players were careful and played only undominated best replies. Then they would finally be playing only (T, L) .

This example shows a drawback of the notion of minimal curb sets: They can contain strategies that are weakly dominated. Therefore let us recall from Basu and Weibull (1991) the notion of sets that are closed under undominated best replies. Formally, s_i is weakly dominated by s'_i if $u_i(s_i, s_{-i}) \leq u_i(s'_i, s_{-i})$ for all s_{-i} with strict inequality for at least one s_{-i} . Let $\text{UBR}(\sigma)$ denote the set of pure best replies against σ that are not weakly dominated.

	L	R
T	1,1	1,1
B	1,1	0,0

FIGURE 4

DEFINITION 5. A non-empty Cartesian set $C = \prod_{i=1}^n C_i$ is closed under undominated best replies (or C is a curb* set) if $\text{UBR}(\prod_{i=1}^n \Delta(C_i)) \subset C$. Such a set is called a minimal curb* set if it does not properly contain a set that is closed under undominated best replies. Strategies contained in minimal curb* sets are called curb* strategies.

LEMMA 2. Every curb set contains a curb* set. Every minimal curb* set contains the support of a Nash equilibrium. Curb* strategies are not weakly dominated.

Proof. Easy and hence omitted. ■

It is easy to adjust the learning process so that players will end up playing curb* strategies. Just replace “best replies” by “undominated best replies” and analogies of Theorems 2, 3, and 4 can be proved easily. On the level of Bayesian players this means that, although they have certain beliefs, they are not completely sure that these beliefs are “correct.”³ Therefore they should be careful and play only undominated best replies.

The approach taken above is a bit unsatisfactory since the uncertainty is not modeled. We will do that now. Remember the sampling procedure described in Section 3. Every time an individual is drawn from class V_i , he hears about L precedents concerning the way player j played this game before. This sample is transformed (by the fictitious play rule) into a belief μ^i from the L -grid distribution space $B^i(h, L)$, where h denotes the recent history of plays.

Now suppose that the final belief of this player is not necessarily μ^i , but some $\hat{\mu}^i$ “close” to μ^i , reflecting the uncertainty of this player. This uncertainty may stem from the fact that the player realizes that he only draws a sample and that μ^i is only a point estimate of the distribution of strategies. The final belief $\hat{\mu}^i$ could be a draw from some “confidence interval” around μ^i . This draw might depend on personal characteristics, as well as on other external factors. We will just assume $\hat{\mu}^i$ is drawn from the uniform distribution over $B_\varepsilon(\mu^i) = \{\sigma^i \in \Sigma_{-i} \mid d_{\max}^i(\mu^i, \sigma^i) \leq \varepsilon\}$, where $\varepsilon > 0$ is fixed⁴ and where $d_{\max}^i(\mu^i, \sigma^i) = \max_{s_{-i} \in S_{-i}} |\mu^i(s_{-i}) - \sigma^i(s_{-i})|$. Note that, for large L , the union of these intervals over all L -grid distributions induced by h consists of all probability distributions close to $\prod_{j \neq i} \Delta(\pi_j(h))$.

What consequences does this have for our learning process? Or, in other words, what strategies will be played with positive probability after each possible history? Well, let $h \in H$ and let $s_i \in S_i$. Before we had that s_i was

³ The uncertainty of the players could stem from the fact that players may realize that other players have different samples. Anyway, sometimes players “are right” to be uncertain since it is possible that a history h is followed by the play of s , where $s \notin \text{span}(h)$.

⁴ We could take $\varepsilon = 1/L$ to reflect the intuition that larger samples should result in smaller confidence intervals.

played with positive probability whenever there was some $\mu^i \in \prod_{j \neq i} \Delta(\pi_j(h))$ such that $s_i \in \text{BR}_i(\mu^i)$. Now we have that s_i is played with positive probability only if the stability region of s_i ,

$$\text{St}_i(s_i) = \{\sigma_{-i} \in \sum_{-i} | s_i \in \text{BR}_i(\sigma_{-i})\},$$

has positive probability under the uniform distribution over $B_\epsilon(\hat{\mu}^i)$, for some L -grid distribution $\hat{\mu}^i$ induced by h . For sufficiently large L , this is equivalent to

$$\mu_{-i} \in \text{cl}(\text{int}(\text{St}_i(s_i))) \tag{5.3}$$

for some $\mu_{-i} \in \prod_{j \neq i} \Delta(\pi_j(h))$, where $\text{cl}(\cdot)$ and $\text{int}(\cdot)$ stand for closure and interior (in the topological space Σ_{-i}), respectively.

Note that if $\mu_{-i} \in \text{int}(\text{St}_i(s_i))$, then s_i is a best reply against each strategy in an open neighborhood of μ_{-i} . Up to equivalence, s_i is then also the unique (and undominated) best reply against this neighborhood, and s_i is called a robust best reply against μ_{-i} . If only (5.3) is satisfied, there is some non-empty open set close to μ_{-i} against which s_i is the unique best reply, and we call s_i a semi-robust best reply against μ_{-i} , which is denoted by $s_i \in \text{SRBR}_i(\mu_{-i})$. As opposed to robust best replies, semi-robust best replies always exist, and there may exist several semi-robust best replies against some μ_{-i} , even if player i has no equivalent strategies. It is easy to see that semi-robust best replies are not weakly dominated. Similar to the case with the (undominated) best reply correspondence we define

DEFINITION 6 (Balkenborg, 1992). A non-empty cartesian set $C = \prod_{i=1}^n C_i$ is closed under semi-robust best replies (or C is a robust set) if $\text{SRBR}(\prod_{i=1}^n \Delta(C_i)) \subset C$. Such a set is called a minimal robust set if it does not properly contain a set that is closed under semi-robust best replies.

It is easy to see that every curb* set contains a robust set, but not every minimal robust set is (contained in) a minimal curb* set. Moreover, every robust set contains the support of a Nash equilibrium.

The learning process where players are uncertain can be described by a Markov chain that is very similar to the ones we had before. Just replace “best replies” by “semi-robust best replies” and analogies of Theorems 2, 3, and 4 can be proved easily. Play will settle down in a minimal robust set.

For “generic” normal form games the minimal curb, curb*, and robust sets coincide with the persistent sets. Persistent sets consist of the extreme points of persistent retracts (Kalai and Samet, 1984). As a matter of fact, for games in which no player has equivalent strategies, the minimal robust sets coincide with the persistent sets (see Balkenborg, 1992). However,

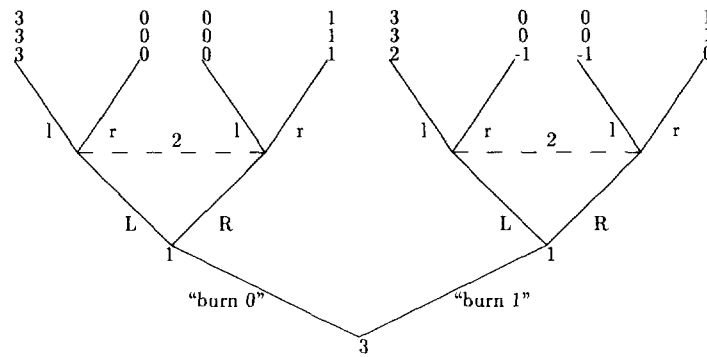


FIGURE 5

many normal form games are interesting because they are the normal form representation of an extensive form game, and these are not “generic” in the class of normal form games. This is due to the fact that there may be strategies in the extensive form game that preclude some information sets (or subgames) from being reached. This implies that curb sets may differ from robust sets. To illustrate this difference consider the following example that is taken from Hurkens (1993).

Consider the game in Fig. 5. Player 3 can decide to burn one unit before players 1 and 2 play a simultaneous move coordination game. Consider the strategy profile $s^{\text{ineff}} = (RR, rr, \text{“burn 0”})$. The singleton set containing this profile is persistent and robust: Player 3 has a unique best reply against s^{ineff} , namely “burn 0”; players 1 and 2 have a lot of (undominated) best replies against s^{ineff} , but they have a unique semi-robust best reply. In a small neighborhood outside $\{s^{\text{ineff}}\}$, players 1 and 2 have a unique best reply, since they have an interest in choosing the same action: in a small neighborhood player 1 plays R with a very high probability, whether or not player 3 burnt something, and hence player 2 has to choose r , whether or not player 3 burnt something.⁵ Since players 1 and 2 have a lot of (undominated) best replies against s^{ineff} , it is easy to see that $\{s^{\text{ineff}}\}$ is not curb or curb*. In fact, the only minimal curb (or curb*) set of this game consists of all strategy profiles yielding the payoff vector (3, 3, 3). The latter set is also persistent and robust.

It seems that the learning processes introduced in this paper are a bit peculiar in the case of extensive form games. In the process that leads to minimal curb sets, sometimes players are absolutely sure that a particular information set will not be reached. Therefore they are free to choose any

⁵ This result depends on the assumption of independent beliefs. See also subsection 5.3.

action in this information set. On the other hand, if we add a little bit of uncertainty, players are still quite certain about the strategies that will be used, but they are also certain that all information sets will be reached with positive probability. Therefore they have to play a best reply against the strategy profile that they believe to be played almost certainly, in all information sets, although many of these information sets will not be reached if this strategy profile is indeed played.

These peculiarities are due to our assumption about the information that players have. In our learning process we assumed that players know the strategies played in the past. For extensive form games it makes more sense to assume that players observe only the *outcomes* of actual play and that they may hold any beliefs about strategies in unreached information sets. We deal with this issue in Section 6.

5.3. *Dependent Beliefs*

Throughout this paper we assumed that a player's belief about the strategies of the other players is independent, i.e., is an element of Σ_{-i} . This was a consequence of the sampling procedure we described in Section 3. Players receive information about the strategies of the players individually. Moreover, if players realize that the players are deciding simultaneously and independently, then it is natural to have independent beliefs. There are, however, two problems concerning the independency of beliefs.

First, do players indeed decide independently? After all, the choices of all players depend (via the samples) indirectly on the same recent history. History might act as a correlation mechanism. Second, our other interpretation of the learning process was that personal characteristics are important to form beliefs. All players expect that strategies that have not been played recently will not be played, but different players may have different assessments of the probabilities with which the remaining strategies are played. In view of this interpretation, an individual from class V_i might have a dependent belief, i.e., an element of $\Delta(S_{-i})$. For instance, he might believe that the other players can correlate their strategies. It does not really matter whether or not the other players do correlate; what matters is that some individuals may believe that they do.

In this section we will examine the consequences of allowing players to have dependent beliefs. We will assume that the classes are very diverse: If h denotes the recent history and $s_i \in BR_i(\mu^i)$ for some $\mu^i \in \Delta(\text{span}(h))$, then s_i will be played with positive probability. Again, we will define a whole set of transition matrices describing such learning processes. Let $\mathcal{B}^{dep}(h) = \{\mu \in \Delta(S) \mid \text{supp}(\mu) \subset \text{span}(h)\}$ denote the set of all dependent beliefs a player may have.

DEFINITION 7. Let \mathcal{P}^{dep} denote the set of transition matrices P that satisfy for all histories $h, \hat{h} \in H$,

$$P(h, \hat{h}) > 0 \Leftrightarrow \begin{cases} \hat{h} \text{ is a successor of } h, \text{ and} \\ r(\hat{h}) \in \text{BR}(\mathcal{B}^{\text{dep}}(h)). \end{cases}$$

Remark. Note that our definition of the transition matrices does not correspond to what one may call “correlated learning.” Suppose that in a three player game player 3 observes that the other players played TL and BR in the last two periods. Then, under our assumption of dependent beliefs, it is possible that player 3 believes that TR and BL will be played, both with probability $\frac{1}{2}$. One may feel that only beliefs of the form $\alpha TL + (1 - \alpha)BR$ should be allowed. We do not know whether such “correlated learning” processes converge to some static set-valued solution concept.

We can prove a theorem similar to Theorem 2. Of course, the process will in general not converge to a minimal curb set, but to a Cartesian set $F = \prod_{i=1}^n F_i$ that is minimal will respect to the following property: If $\mu \in \Delta(F)$ and $s_i \in \text{BR}_i(\mu_{-i})$, then $s_i \in F_i$. Following Harsanyi and Selten (1988) we call such a set a primitive formation.⁶

THEOREM 5. *There exists $\underline{K} \in \mathbb{N}$ such that for all $K \geq \underline{K}$ and for every Markov chain with a transition matrix $P \in \mathcal{P}^{\text{dep}}$:*

- (i) *If $Z \subset H$ is an ergodic set then $Z \subset F^K$ for some primitive formation F .*
- (ii) *For every primitive formation F there exists exactly one ergodic subset $Z \subset F^K$.*

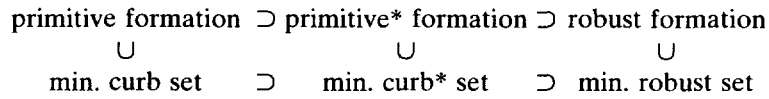
We omit the proof because it is essentially the same as the proof of Theorem 2. We just have to observe that if F is a primitive formation and $s \in F$, then s_j is a best reply against some (dependent) belief concentrated on F .

Obviously, analogies of Theorems 3 and 4 to the case of dependent beliefs also exist. The same is true for the results of Section 5.2 on undominated best replies and semi-robust best replies. Analogous to curb* and robust sets we could define primitive* and robust formations. The reader should be aware, though, that the definition of semi-robustness needs to be adapted. In the context of dependent beliefs we say that s_i is a semi-robust best reply against $\mu_{-i} \in \Delta(S_{-i})$ if $\mu_{-i} \in \text{cl}(\text{int}(\text{St}_i(s_i)))$, where $\text{cl}(\cdot)$ and $\text{int}(\cdot)$ stand for closure and interior, respectively, in the topological space $\Delta(S_{-i})$, and where $\text{St}_i(s_i) = \{\mu_{-i} \in \Delta(S_{-i}) \mid s_i \in \text{BR}_i(\mu_{-i})\}$ is the stability region of s_i .

Of course, in a two-person game the primitive formations are identical

⁶ Harsanyi and Selten (1988) consider this concept in the agent normal form.

to the minimal curb sets. Moreover, every primitive formation contains a minimal curb set. Hence, if a game has a unique minimal curb set C which is also a primitive formation, then C is also the unique primitive formation.⁷ Similar statements can be made about the other concepts with the help of the following diagram. In this diagram, $X \supset Y$ means that every X contains an Y , but not every Y is contained in an X .



6. LEARNING FROM OUTCOMES

Throughout this paper we assumed that players know the strategies that were used in the past. This assumption is reasonable when the players in the underlying game choose their actions simultaneously. But if the underlying game is in fact an extensive form game, it makes more sense to assume that players observe only the outcomes, i.e., the paths in the tree generated by the strategies. Consider for example the “burning money” game in Fig. 5. Suppose player 3 chose to “burn 0” in the last period. How could he know how players 1 and 2 would have reacted to “burn 1”? In fact, he cannot, although he may have some beliefs.

In this section we will consider the case where players observe only the outcomes in the recent past. We assume that all agents form expectations on the basis of observed outcomes, and that different agents within a pool may form different beliefs. We pose only one restriction on the beliefs: When a player is able to conclude from the observed outcomes that a particular strategy has not been played during the last K periods, he expects it will not be played next period. As before, we assume that the classes are very diverse: As soon as strategy s_i is a best reply against some independent belief, satisfying this restriction, s_i will be played with positive probability.

We will define a class of transition matrices that correspond to such a “learning from outcomes” process, and we denote this class by \mathcal{P}^{out} . Before we can do so, we need some notation.

Let G be an extensive form game. Let \mathcal{O} denote the set of outcomes (i.e., paths in the tree from the root to an endpoint) and let $o: S \rightarrow \mathcal{O}$ be the mapping that assigns to a pure strategy combination the outcome it generates. We will assume that there are no moves of nature in G , since this

⁷ For an example, consider the game shown in Fig. 5. The set of strategy profiles yielding the payoff vector (3, 3, 3) is the unique minimal curb set, but it is also a primitive formation. Moreover, it is also the unique primitive* and robust formation.

mapping is not well defined if there are. For a history $h = (s^{-K}, \dots, s^{-1})$, let $\text{outc}(h) = \{o(s^{-K}), \dots, o(s^{-1})\}$. Note that $\text{outc}(h)$ summarizes the information a player has. Let $\text{cons}_i(h) = \{s_i \in S_i \mid \exists s_{-i} \in S_{-i} \text{ s.t. } o((s_i, s_{-i})) \in \text{outc}(h)\}$ denote the set of strategies of player i that are consistent with the observed outcomes. Let $\text{cons}(h) = \prod_{i=1}^n \text{cons}_i(h)$.

DEFINITION 8. Let \mathcal{P}^{out} denote the set of transition matrices P that satisfy for all histories $h, \hat{h} \in H$,

$$P(h, \hat{h}) > 0 \Leftrightarrow \begin{cases} \hat{h} \text{ is a successor of } h, \text{ and} \\ r(\hat{h}) \in \text{BR}(\mathcal{B}^{\text{ind}}(\text{cons}(h))). \end{cases}$$

In general, it is not true that play will settle down in minimal curb sets. Note that $\text{cons}(h) \supset \text{span}(h)$. This implies that if $P \in \mathcal{P}$ and $P(h, \hat{h}) > 0$, then $P^{\text{out}}(h, \hat{h}) > 0$ for all $P^{\text{out}} \in \mathcal{P}^{\text{out}}$. Using part of the proof of Theorem 2, it follows that, if K is large enough, for every history h and every $P^{\text{out}} \in \mathcal{P}^{\text{out}}$, there exists a curb history \tilde{h} such that $h \rightsquigarrow \tilde{h}$. The problem is that there might exist a history \hat{h} , which is not a curb history, such that $\tilde{h} \rightsquigarrow \hat{h}$. This might even happen in “generic” extensive form games, as the game presented in Fig. 6 shows.

This game has a unique minimal curb set, namely $\{U, D\} \times \{aA, aB, aC, bA\}$. However, suppose that in the recent (curb) history the strategy combinations (D, aB) and (U, bA) were played. Hence, player 1 observes (among other things) the outcomes DB and Ub . He might believe that the strategy bB was played and will be played again next period. If he does so, he will choose “Out,” which is not a curb strategy.

The above example seems to suggest that there is no hope to obtain a result like Theorem 2 in the case of learning from outcomes. There are, however, two classes of games for which such an analogy does exist. The

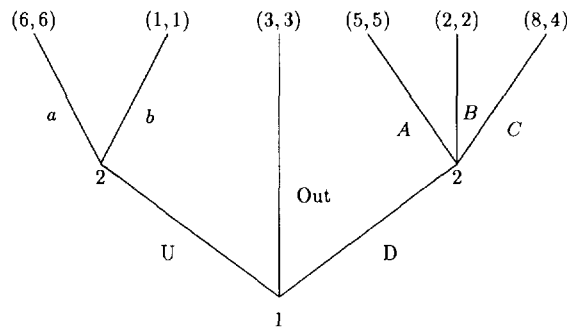


FIGURE 6

first class consists of the extensive form games without moves of nature, where each player has only one information set at which he has to make a choice. For obvious reasons we call such a game an agent normal form game without moves of nature, and we denote the class by ANF. The second class of games consists of those games G that have the property that any minimal curb set C of G corresponds to a single outcome; i.e., the set $\{o(c) | c \in C\}$ is a singleton. We denote this class by SCO (single curb outcome). Examples of these games are shown in Figs. 1b, 1c, and 5.

To prove the above claims we just need to show that $\mathcal{P}^{\text{out}} \subset \tilde{\mathcal{P}}$, where $\tilde{\mathcal{P}}$ is as defined at the beginning of Section 5. Part (5.1) follows from $\text{span}(h) \subset \text{cons}(h)$ and part (5.2) follows from the next lemma.

LEMMA 3. *Let $G \in \text{ANF}$ or $G \in \text{SCO}$ and let C be a minimal curb set of G . Then*

$$h \in C^K \Rightarrow \text{cons}(h) \subset C.$$

Proof. First consider the case $G \in \text{ANF}$. Let j be a player. If there is an outcome $o(s^{-m}) \in \text{outc}(h)$ that does not intersect j 's information set, then it follows that $\text{BR}_j(s^{-m}) = S_j$. This implies that $C_j = S_j \supset \text{cons}_j(h)$. If there is no such outcome, all outcomes intersect j 's information set and $\text{cons}_j(h) = \pi_j(h) \subset C_j$. Hence, $\text{cons}(h) \subset C$.

Now consider the case $G \in \text{SCO}$. Let $\bar{s} = r(h)$. Now we have $\text{outc}(h) = \{o(\bar{s})\}$. Let j be a player and suppose $s_j \in \text{cons}_j(h)$. In any information set of j that intersects $o(\bar{s})$, s_j picks the same action as \bar{s}_j , since s_j is consistent with h . Since $G \in \text{SCO}$, we have that \bar{s}_j is a best reply against \bar{s}_{-j} . But this implies that s_j is a best reply against \bar{s}_{-j} as well, and hence $s_j \in C_j$.

The reader can check that there are also analogies of Theorems 3 and 4 to the case where players learn from outcomes. The definition of a mimicker needs to be adapted, since players do not observe strategies. We may assume that mimickers choose at random a strategy from the set of strategies that are consistent with (some of) the observed outcomes. There is also an analogy of Theorem 5, where players' beliefs are not independent. There are, however, no analogues for the results of Section 5.2 on the refined notions of undominated best replies or of semi-robust best replies. This is due to the fact that strategies that are consistent with a curb* history may be weakly dominated. The game presented in Fig. 6 shows an example of such a case: The only curb* strategy is (U, aA) , but aB and aC are consistent with the curb* outcome.

7. LEARNING AND EXPERIMENTATION

In many papers on learning, experimentation plays a prominent role. (See, e.g., Kandori *et al.*, 1993; Samuelson, 1994; Young, 1993; Fudenberg and Kreps, 1988).

In Young (1993), Samuelson (1994), and Kandori *et al.* (1993) the possibility of experimentation (or mistakes or mutations) implies that the Markov chain describing the learning process becomes irreducible and hence has a unique stationary distribution. By taking the limit as the experimentation rate tends to zero, one stationary distribution of the unperturbed process is selected. In Young (1993) and Kandori *et al.* (1993) this yields typically a unique so called stochastically stable state because they consider a special class of games. Samuelson (1994) considers games with alternative best replies and then the support of the limit distribution consists usually of one or more line segments.

It turns out that the introduction of experimentation does not change the results of the present paper, at least not for two-person games. If a two-person game has multiple minimal curb sets, experimentation will not yield the selection of a particular one: the limiting distribution puts positive weight on all states that are ergodic under the unperturbed process. The intuition behind this result is that only one mistake by one player is necessary in order to move the system from one ergodic set to another. When the game has more than two players, it might happen that a particular minimal curb set is selected. One can characterize the selected minimal curb set graph-theoretically.

In order to prove these results formally, we would have to recall the essential definitions and theorems from Young (1993). We refer the reader to the original paper for a formal treatment. We will just illustrate the result by means of an example.

Consider again the coordination game from Fig. 1a. As we have seen before, this game has two minimal curb sets, $\{(T, L)\}$ and $\{(B, R)\}$. Suppose the system is in state $h^{TL} = (TL, \dots, TL)$ and player 1 makes a mistake and plays B . Since sampling occurs with replacement, player 2 may receive a draw containing many B 's, in which case he will play R . It may happen that from then on player 1 receives draws with many R 's while player 2 keeps drawing many B 's. It follows that, after the initial mistake, the system can move to $h^{BR} = (BR, \dots, BR)$, without making any further mistakes. Hence, only one mistake is needed to move the system from h^{TL} to h^{BR} . Similarly, only one mistake is needed to move the system from h^{BR} to h^{TL} . Since the mistake probabilities are of the same order, the limiting distribution puts positive weight on both ergodic states.

This result is in contrast with Young (1993). In Young (1993) the players also have information about play in the recent history: Every player draws

a sample of m plays out of the plays of the most recent K periods, but without replacement. Then players play a best reply in a fictitious play fashion. Consider again the coordination game from Fig. 1a. Suppose that the system is in state h^{TL} and that player 1 makes a mistake and plays B . If no further mistakes occur, the system will move back to h^{TL} if the sample size is at least 2: Since sampling occurs without replacement, every sample contains at least as many T 's as B 's, and player 2 will always play L (unless he makes a mistake). It is easy to see that in this example at least $3m/4$ mistakes are needed to move the system from h^{TL} to h^{BR} , while only $m/4$ mistakes are needed to move the system in the other direction. It follows that h^{TL} is the unique stochastically stable state.

8. CONCLUDING REMARKS

We have considered learning processes where the players have a bounded memory and play best replies against past play. The importance of the bounded memory can be elucidated by comparing our learning process with Milgrom and Roberts (1991). In general they consider games with compact strategy sets that are played continuously. Translated to the context of a two-player finite normal form game which is played repeatedly at discrete points in time, they define a sequence of plays $\{s(t)\}_{t=0}^{\infty}$ to be *consistent with adaptive learning* if for all \hat{t} there exists a \bar{t} such that for all $t \geq \bar{t}$, $s(t+1) \in \text{BR}(\mathcal{B}^{\text{ind}}(\{s(\hat{t}), s(\hat{t}+1), \dots, s(t)\}))$. We could similarly define this sequence to be *consistent with learning with bounded memory* if there exists $K \in \mathbb{N}$ such that for all t , $s(t+K) \in \text{BR}(\mathcal{B}^{\text{ind}}(\{s(t), s(t+1), \dots, s(t+K-1)\}))$. This definition illustrates the similarity between the present paper and Milgrom and Roberts (1991).

Consider for example the pure coordination game shown in Fig. 1a. The sequence $TR, BL, TR, BL, TR, \dots$ satisfies both definitions of consistency. However, the finiteness of the memory and of the strategy space allows us to obtain a finite Markov chain, from which we can compute that the probability of obtaining the above sequence is zero: Only sequences with tails TL, TL, TL, \dots or BR, BR, BR, \dots are obtained with positive probability.

Milgrom and Roberts (1991) show that sequences that are consistent with adaptive learning will eventually lie within the set of serially undominated strategies, which is a superset of the set of rationalizable strategies. They give some examples of games with strategic complementarities where this set is a singleton, which implies that these sequences must converge to the unique equilibrium. We get the same results in these games because the set of curb strategies is a subset of the set of rationalizable strategies. But we get similar results in some games where the set of rationalizable strategies is

large. In every game that has a unique and strict equilibrium \bar{s} , $\{\bar{s}\}$ is the unique minimal curb set. Hence, in such games our learning process leads the players (with probability 1) to the unique equilibrium (Corollary 2). An example of such a game is given in Fig. 3, where all strategies are rationalizable.

Another example is the discretized version of the following three player Cournot oligopoly game. Player i chooses to produce q_i at zero costs to maximize $q_i(A - q_1 - q_2 - q_3)$. The unique (and strict) equilibrium is $(A/4, A/4, A/4)$. The set of rationalizable strategies is $[0, A/2] \times [0, A/2] \times [0, A/2]$.

REFERENCES

- BALKENBORG, D. (1992). *The Properties of Persistent Retracts and Related Concepts*, Ph.D. thesis, University of Bonn.
- BALKENBORG, D. (1993). "Strictness, Evolutionary Stability and Repeated Games with Common Interests," CARESS Working Paper 93-20.
- BASU, K., AND WEIBULL, J. W. (1991). "Strategy Subsets Closed under Rational Behavior," *Econ. Lett.* **36**, 141-146.
- BERNHEIM, B. D. (1984). "Rationalizable Strategic Behavior," *Econometrica* **52**, 1007-1028.
- BLUME, A. (1993a). "Neighborhood Stability in Sender-Receiver Games," *Games Econ. Behav.*, in press.
- BLUME, A. (1993b). "Communication, Risk and Efficiency in Games," mimeo.
- BROWN, G. W. (1951). "Iterative Solution of Games by Fictitious Play," in *Activity Analysis of Production and Allocation*. New York: Wiley.
- VAN DAMME, E., AND HURKENS, S. (1993). "Commitment Robust Equilibria and Endogenous Timing," CentER discussion paper no. 9356.
- FUDENBERG, D., AND KREPS, D. (1988). "A Theory of Learning, Experimentation, and Equilibrium in Games," mimeo, Stanford University and Massachusetts Institute of Technology.
- HARSANYI, J., AND SELTEN, R. (1988). *A General Theory of Equilibrium Selection in Games*. Cambridge, MA: MIT Press.
- HURKENS, S. (1993). "Multi-sided Pre-play Communication by Burning Money," *J. Econ. Theory*, in press.
- KALAI, E., AND SAMET, D. (1984). "Persistent Equilibria in Strategic Games," *Int. J. Game Theory* **13**, 129-144.
- KALAI, E., AND SAMET, D. (1985). "Unanimity Games and Pareto Optimality," *Int. J. Game Theory* **14**, 41-50.
- KANDORI, M., MAILATH, G. J., AND ROB, R. (1993). "Learning, Mutation, and Long Run Equilibria in Games," *Econometrica* **61**, 29-56.
- KEMENY, J., AND SNELL, J. (1976). *Finite Markov Chains*. New York/Heidelberg/Berlin: Springer-Verlag.
- MILGROM, P., AND ROBERTS, J. (1991). "Adaptive and Sophisticated Learning in Normal Form Games," *Games Econ. Behav.* **3**, 82-100.

- PEARCE, D. G. (1984). "Rationalizable Strategic Behavior and the Problem of Perfection," *Econometrica* **52**, 1029–1050.
- ROBINSON, J. (1951). "An Iterative Method of Solving a Game," *Ann. Math.* **54**, 296–301.
- SAMUELSON, L. (1994). "Stochastic Stability in Games with Alternative Best Replies," *J. Econ. Theory* **64**, 35–65.
- YOUNG, H. P. (1993). "The Evolution of Conventions," *Econometrica* **61**, 57–84.