

Implementation, Elimination of Weakly Dominated Strategies and Evolutionary Dynamics¹

Antonio Cabrales

*Departament d'Economia i Empresa, Universitat Pompeu Fabra, Ramon Trias Fargas,
25-27, 08005 Barcelona, Spain; and
ELSE, Centre for Economic Learning and Social Evolution, University College London,
Gower Street, London WC1E 6BT, United Kingdom*

and

Giovanni Ponti²

*Departamento de Fundamentos del Análisis Económico, Universidad de Alicante, 03071
Alicante, Spain; and
ELSE, Centre for Economic Learning and Social Evolution, University College London,
Gower Street, London WC1E 6BT, United Kingdom*

Received August 11, 1997

This paper studies convergence and stability properties of T. Sjöström's (1994, *Games Econom. Behav.* 6, 502–511) mechanism, under the assumption that boundedly rational players find their way to equilibrium using monotonic evolutionary dynamics and best-reply dynamics. This mechanism implements most social choice functions in economic environments using as a solution concept one round of deletion of weakly dominated strategies and one round of deletion of strictly dominated strategies. However, there are other sets of Nash equilibria, whose payoffs may be very different from those desired by the social choice function. With monotonic dynamics, all these sets of equilibria contain limit points of the evolutionary dynamics. Furthermore, even if the dynamics converge to the “right”

¹ The authors are grateful to Ted Bergstrom, Ken Binmore, Sandro Brusco, Tilman Börgers, Leo Hurwicz, and Karl Schlag for stimulating comments, as well as suggestions by an anonymous referee on earlier versions of this paper. The usual disclaimers apply. Both authors gratefully acknowledge financial support from the European Commission, under the “Human Capital Mobility” program. Financial support from the Spain's Ministry of Education, under Grant PB96-0302 (Cabrales) and from the Instituto Valenciano de Investigaciones Economicas (Ponti) is also gratefully acknowledged. This paper was written while Cabrales was visiting the Centre for Economic Learning and Social Evolution (ELSE) at the University College London. He would like to thank ELSE for hospitality and encouragement.

² Corresponding author. Departamento de Fundamentos del Análisis Económico, Universidad de Alicante, 03071 Alicante, Spain. E-mail: giuba@merlin.fae.ua.es.

set of equilibria (i.e., the one which contains the solution of the mechanism), it may converge to an equilibrium which is worse in welfare terms. In contrast with this result, any interior solution of the best-reply dynamics converges to the equilibrium whose outcome the planner desires. *Journal of Economic Literature* Classification Numbers: C72, D70, D78. © 2000 Academic Press

Key Words: implementation theory, evolutionary dynamics, bounded rationality.

1. INTRODUCTION

The theory of implementation studies the problem of designing decentralized institutions (“mechanisms”) through which certain socially desirable objectives can be achieved. More precisely, a *social choice rule* is implemented by a (*game-form*) mechanism if, for every possible environment (preference profile), the solution (set of equilibrium outcomes) of the mechanism coincides with the set of outcomes of the social choice rule.

This definition implicitly assumes that agents are always able to play equilibrium strategies. However, there is substantial empirical and experimental evidence against this theoretical presumption.³ In spite of this evidence, research in implementation theory has paid little attention to the problem of how equilibrium is achieved.⁴ Since the planner should be concerned with the performance of the mechanism when some (if not all of the) agents are not as “rational” as expected, it is useful to test the mechanism’s performance in the presence of some form of bounded rationality.

A more fundamental approach to these issues would require the planner to take bounded rationality into account, when designing the game agents play. This necessarily leads to an alternative definition of implementation which includes, among the variables which specify the “environment,” the learning protocols agents use, as well as initial conditions of the dynamic process. In this respect, we propose the following definition. For a given set of environments Φ and a given set of dynamics D , a social choice rule is *dynamically implemented* by a mechanism if, for all $\phi \in \Phi$ and $d \in D$, the limiting set of outcomes coincides with the set of outcomes of the social choice rule.

There is a caveat here. Why should we focus only on limiting outcomes? The planner may also care about what happens on the way to equilibrium, as the dynamic path may include outcomes significantly different than

³ See Cooper et al. (1991) for the prisoner’s dilemma, a strictly dominance solvable game, McKelvey and Palfrey (1991) for the centipede game, a game with a unique Nash equilibrium, and Güth et al. (1982) for the ultimatum game, which has a unique subgame perfect equilibrium.

⁴ Noticeable exceptions are the papers of Muench and Walker (1984), Walker (1984), Jordan (1986), Vega-Redondo (1989), De Trenqualye (1988, 1989), and Cabrales (1999).

what the choice rule prescribes. This, in turn, would require one to fully characterize the planner's preferences, rather than specify the most preferred outcome, for any given state of the environment. This is something the implementation literature traditionally leaves unspecified. Moreover, if the planner does not discount the future and the game is played infinitely often, then it is legitimate to look at limiting outcomes.

In this paper we study the *dynamic* implementation of Sjöström's (1994) mechanism.⁵ We concentrate on Sjöström's mechanism for several reasons. First, the conditions for implementation are quite weak. Although the environments that are permitted are not universal, they are rich enough for most economic purposes. Furthermore, this reduction in the domain allows the author to implement the social choice rule with a "bounded" game, that is, a game which does not exploit equilibrium nonexistence to rule out undesirable outcomes.⁶ Finally, the game can be solved by one round of deletion of weakly dominated strategies, and then another round of deletion of strictly dominated strategies. This feature of the mechanism makes it particularly attractive since, under some assumptions of imperfect knowledge of agents,⁷ the appropriate solution concept implies one round of deletion of weakly dominated strategies, and then the iterated deletion of strictly dominated strategies.

In Sjöström's (1994) mechanism agents are arranged to simultaneously announce their own preferences, together with the preferences of their two closest neighbors. The mechanism is designed in such a way that the truthful report of one's own preferences is weakly dominant, as it does not affect one's payoff, except for a set of (so-called) *totally inconsistent* states, where it is (strictly) preferable to report preferences truthfully. Since, for this mechanism, it is always advantageous to report the same preferences about your neighbors as what they are reporting about themselves, it is clear that the only equilibrium that survives the first round of deletion of weakly dominated strategies is the truth-telling one.

However, there are many other Nash equilibria. In particular, for every preference profile R , there is a component (i.e., a closed and connected set) of equilibria in which all agents report the preferences for their neighbors indicated in R , and report the preferences about themselves indicated in R with high enough (this need not be very high) probability.

⁵ Sjöström's (1994) mechanism and the one proposed by Jackson et al. (1994) for separable environments are very similar. Most of our results would generalize easily for that mechanism as well.

⁶ For example, in the canonical mechanism for Nash implementation (Repullo, 1987), if agents disagree widely on the announced preferences, they have to play a game in which the agent announcing the highest integer wins a prize.

⁷ Either because of payoff uncertainty, as in Dekel and Fudenberg (1990), or through lack of common knowledge of rationality, as in Börgers (1994).

This is because it is important for the mechanism that all agents match their neighbors' announcements about themselves, but the report about oneself is only important in some unlikely (totally inconsistent) state.

First, we study the performance of the mechanism under monotonic dynamics (Nachbar, 1990), which essentially imply higher growth rates for those strategies which perform better.⁸ For these dynamics, we show (Proposition 4) that many equilibria in all equilibrium components are limit points of trajectories of the evolutionary dynamics that have completely mixed initial conditions (that is, initial conditions that give strictly positive weights to all possible messages). Even when the dynamics converge to the "right" component of equilibria (i.e., the one which contains the solution of the mechanism), they need not go to the "right" equilibrium. We also show by example (Proposition 2) that the initial conditions that lead to these equilibria need not be close to the limiting point. We also study how the dynamic structure reacts to the introduction of (arbitrarily small) perturbations in the vectorfield. In the example we show (Proposition 6) that, although there is a unique structurally stable component (namely, the component which contains the undominated equilibrium of the game), the untruthful component is stable for a non-negligible set of admissible perturbations.

In other words, the less responsive the dynamics are to payoffs (the further the initial conditions from the "right" equilibrium), the more difficult it is to converge to the desired solution. Only in the extreme case of best-reply dynamics (in which the response to arbitrarily small payoff differences is infinite), we show (Proposition 7) that any interior trajectory converges to the pure strategy equilibrium in which players reveal their true preferences and the outcome desired by the planner is achieved.

The fact that evolution need not eliminate weakly dominated strategies has been known since, at least, Nachbar (1990). However, we are far from possessing a sound theory on the evolutionary properties of weakly dominance solvable games, as we have examples in which a single round of deletion is not allowed if we want to characterize the limiting set of the evolutionary dynamics (see, for example, Samuelson, 1993 and Cressman and Schlag, 1998), as well as games in which only strategies which survive an (arbitrarily large) number of rounds of deletion can be in the support of the limiting play (see, for the finitely repeated prisoners' dilemma, Cressman, 1996; or for the centipede game, Ponti, 2000). Since the theory has not proposed, so far, a suitable framework to explain these differences, it is

⁸ One particularly well known member of the family of monotonic dynamics is the so-called *replicator dynamics* of evolutionary game theory (Taylor and Jonker, 1978). These dynamics have been given a learning theoretic foundation by Börgers and Sarin (1997), and they can also be interpreted as a model of imitation (Schlag, 1994).

important to test the evolutionary properties of (game-form) mechanisms in which the iterated deletion of dominated strategies plays such a crucial role. In this respect, our findings are very similar to those of Gale et al. (1995). They analyze the classic *chain store* game, another game which has a Nash equilibrium component in which a player selects a weakly dominated strategy with positive probability. In both cases, these components are reachable by the evolutionary dynamics, and therefore should not be discarded as a reasonable predictor of the asymptotic play.

The remainder of the paper is arranged as follows. In Section 2 we introduce some notation, we describe the mechanism, and we make the assumptions about the dynamics. In Section 3 we fully characterize (for all interior initial conditions) the set of limit points of any monotonic dynamic for the game in Fig. 1 (Sjöström, 1994) to be considered as a simplified version of the mechanism. In Section 4 we give local results on the convergence and stability properties of the Nash equilibrium components of the general game. In Section 5 we describe the structural stability properties of the equilibria of the simplified mechanism. Section 6 explores the dynamic implementation of Sjöström's (1994) mechanism under best-reply dynamics. Finally, Section 7 concludes, together with an appendix containing the proofs of the relevant propositions.

2. THE MODEL AND THE DYNAMICS

We introduce a few changes to Sjöström's (1994) model for analytical convenience. First, we employ a Von Neumann–Morgenstern utility function instead of a preference relation. This is because we need to specify the payoff functions for mixed strategies, as the dynamics are defined on the mixed strategy space. We also assume that the set of possible preference parameters is finite. This is because the dimension of the pure strategy space is related with the set of preferences. If we had an infinite dimensional pure strategy space, the dynamics, which account for the relative frequency with which each pure strategy is being used, would have to describe the evolution of a measure over an infinite space. This seems an unnecessary complication for our purposes.

There is a set $I \equiv \{1, \dots, n\}$, $n \geq 3$, of agents and a set $A \subseteq \Re_+^m$ of feasible consumption plans. The preferences of agent $i \in I$ are represented with a (Von Neumann–Morgenstern) utility function $v_i: A \times \Phi_i \rightarrow \Re$, where Φ_i specifies a finite set of possible preference parameters. An element R_i of Φ_i represents the preferences of agent i over A . A *preference profile* is a vector $R = (R_1, \dots, R_n)$, which is assumed to be common knowledge among the agents. The following assumptions refine the sets of feasible consumption plans and preference profiles.

Assumption p.1 (free disposal). If $a \in A$ and $0 \leq a' \leq a$, then $a' \in A$.

Assumption p.2. The set of feasible consumption plans A is convex. For all $a, a' \in A$ and for all $\lambda \in [0, 1]$ then $\lambda a + (1 - \lambda)a' \in A$.

Assumption p.3. The preferences represented by $R_i \in \Phi_i$ are strictly convex. For any $a, a' \in R_+^m$ and for all $\lambda \in (0, 1)$, if $a \neq a'$ and $v_i(a, R_i) \geq v_i(a', R_i)$, then $v_i(\lambda a + (1 - \lambda)a', R_i) > v_i(a', R_i)$.

Assumption p.4. For any $R_i \in \Phi_i$ if $a \geq 0$ and $a \neq 0$ then $v_i(a, R_i) > v_i(0, R_i)$.

Assumption p.5 (preference reversal). For any $R_i, R'_i \in \Phi_i$ if $R_i \neq R'_i$ then there are $a, \tilde{a} \in A$ such that $v_i(a, R_i) > v_i(\tilde{a}, R_i)$ and $v_i(\tilde{a}, R'_i) > v_i(a, R'_i)$.

For any set $B \subseteq \mathfrak{R}_+^m$ and any $R_i \in \Phi_i$ a *choice representation* is defined as follows: $c(B, R_i) \equiv \{a \in B \mid \text{for all } b \in B, v_i(a, R_i) \geq v_i(b, R_i)\}$.

For any $i \in I$, a *social choice function* for player i is a mapping $f_i: \Phi \rightarrow A$, where $f(R) \equiv (f_1(R), \dots, f_n(R))$.

Assumption p.6 (individual rationality). For all i and R , $f_i(R) \neq (0, \dots, 0)$.

A *mechanism* is a pair $\Gamma \equiv (M, \alpha)$, where $M \equiv \times_{i \in I} M_i$ and $\alpha: M \rightarrow A$. M_i is the *message space* of agent i (with generic element m_i , and $m = (m_1, m_2, \dots, m_n)$) and α is the *outcome function*. A pair (Γ, R) (a *mechanism* and a *preference profile*) defines a game.

Let $M_{-i} \equiv M_1 \times \dots \times M_{i-1} \times M_{i+1} \times \dots \times M_n$ (with generic element m_{-i}). Given a *mechanism* Γ and a *preference profile* R , we say that m_i is *weakly dominated* for some set of messages $F \equiv \times_{i \in I} F_i \subseteq M$ if there exists a message $m'_i \in F_i$ such that $v_i(\alpha_i(m'_i, m_{-i}), R_i) \geq v_i(\alpha_i(m_i, m_{-i}), R_i)$ for all $m_{-i} \in F_{-i}$ and there is some $m_{-i}^* \in F_{-i}$ such that $v_i(\alpha_i(m'_i, m_{-i}^*), R_i) > v_i(\alpha_i(m_i, m_{-i}^*), R_i)$. Define the set $U_i(F : (\Gamma, R)) \equiv \{m_i \in F_i \mid m_i \text{ is not weakly dominated in } F \text{ for the game } (\Gamma, R)\}$.

The message m_i is a *best response* for player i to $m_{-i} \in M_{-i}$ in the game (Γ, R) , if $v_i(\alpha_i(m_i, m_{-i}), R_i) \geq v_i(\alpha_i(m'_i, m_{-i}), R_i)$ for all $m'_i \in M_i$. A message profile m is a *Nash equilibrium* (NE) for the game (Γ, R) , if m_i is a best response to m_{-i} in the game (Γ, R) for all $i \in I$. A message profile $m \in M$ is an *undominated Nash equilibrium* (UNE) for the game (Γ, R) if it is a Nash equilibrium and $m_i \in U_i(M : (\Gamma, R))$. Let $\text{UNE}(\Gamma, R) \equiv \{\alpha(m) \in A \mid m \text{ is an UNE for the game } (\Gamma, R)\}$.

We say that a mechanism Γ *implements* a social choice function f in *undominated Nash equilibrium* if for all $R \in \Phi$, $f(R) = \text{UNE}(\Gamma, R)$.

For the *iterated deletion of weakly dominated strategies* let $U_i^1(\Gamma, R) = U_i(M : (\Gamma, R))$, and if $U_i^k(\Gamma, R)$ has been defined for $k \geq 1$, let $U_i^{k+1}(\Gamma, R)$

$\equiv U_i(\times_{j \in I} U_j^k(\Gamma, R) : (\Gamma, R))$. Let $U_i^\infty(\Gamma, R) \equiv \bigcap_{k=1}^\infty U_i^k(\Gamma, R)$. Let $\text{IWD}(\Gamma, R) \equiv \{\alpha(m) \in A \mid m_i \in U_i^\infty(\Gamma, R) \text{ for all } i\}$.

We say that a mechanism Γ *implements* a social choice function f with *iterated deletion of weakly dominated strategies* if for all $R \in \Phi$, $f(R) = \text{IWD}(\Gamma, R)$.

We now construct a mechanism.

Let $M_i = \Phi_{i-1} \times \Phi_i \times \Phi_{i+1}$, so that each individual announces the preferences of her two neighbors, and let members of M_i and M be denoted m_i and m , respectively. A generic strategy is therefore $m_i = (R_{i-1}^i, R_i^i, R_{i+1}^i)$. A K -tuple of messages $\{m_{j_1}, \dots, m_{j_K}\}$ is *totally consistent* if, whenever agents $i, k \in \{j_1, \dots, j_K\}$ both announce the preference of player $j \in I$, then $R_j^i = R_j^k$. On the other hand, a K -tuple of messages $\{m_{j_1}, \dots, m_{j_K}\}$ is *totally inconsistent* if, whenever agents $i, k \in \{j_1, \dots, j_K\}$ both announce the preference of player $j \in I$, then $R_j^i \neq R_j^k$.

Consider $R_i, R'_i \in \Phi_i$, where $R_i \neq R'_i$. By Assumption p.5 there are $a, \tilde{a} \in A$ such that $v_i(a, R_i) > v_i(\tilde{a}, R_i)$ and $v_i(\tilde{a}, R'_i) > v_i(a, R'_i)$. We can choose a and \tilde{a} so that $v_i(a, R_i) > v_i(a', R_i)$ for all a' in the line segment between a and \tilde{a} . Given this pair (a, \tilde{a}) let $\beta_i(R_i, R'_i) \equiv \{b \in \mathfrak{R}_+^m \mid b = \lambda a + (1 - \lambda)\tilde{a}, \text{ for } \lambda \in [0, 1]\}$. By construction, for all $R_i, R'_i \in \Phi_i$, $c(\beta_i(R_i, R'_i), R_i) \neq c(\beta_i(R_i, R'_i), R'_i)$. Let $\phi(i, m) \equiv (R_1^n, R_2^n, \dots, R_{i-1}^{i-1}, R_{i+1}^i, R_{i+2}^{i+1}, \dots, R_n^{n-1})$ and, for every i and m_{-i} , define

$$B_i(m_{-i}) = \begin{cases} f_i(\phi(i, m)) & \text{if } m_{-i} \text{ is totally consistent} \\ \beta_i(R_{i-1}^{i-1}, R_{i+1}^{i+1}) & \text{if } m_{-i} \text{ is totally inconsistent} \\ \frac{1}{n} f_i(\phi(i, m)) & \text{otherwise.} \end{cases}$$

Now we can define α :

$$\alpha_i(m) = \begin{cases} c(B_i(m_{-i}), R_i^i) & \text{if } R_{i-1}^i = R_{i-1}^{i-1} \quad \text{and} \quad R_{i+1}^i = R_{i+1}^{i+1} \\ 0 & \text{otherwise.} \end{cases}$$

Let \hat{R} be the true preference profile and let R^* be an arbitrary preference profile. To understand how the mechanism works, notice that the only time when the choice of an announcement R_i^i has any effect on i 's payoffs is when m_{-i} is totally inconsistent. In this case, the outcome is the optimal choice within the set $\beta_i(R_{i-1}^{i-1}, R_{i+1}^{i+1})$ according to the announced R_i^i . This is the reason why, for player i , announcing her true preference \hat{R}_i can never hurt. Furthermore, for every alternative announcement $R_i^i = R_i^*$, there is some totally inconsistent m_{-i} with $R_{i-1}^{i-1} = \hat{R}_i$ and $R_{i+1}^{i+1} = R_i^*$ and the set $\beta_i(., .)$ is constructed in such a way that $c(\beta_i(\hat{R}_i, R_i^*), \hat{R}_i)$ is strictly

preferred to $c(\beta_i(\hat{R}_i, R_i^*), R_i^*)$. Therefore, a message $m_i = (R_{i-1}^i, R_i^*, R_{i+1}^i)$ is weakly dominated by a message $m_i = (R_{i-1}^i, \hat{R}_i, R_{i+1}^i)$; i.e., untruthful announcements about oneself are weakly dominated.

Once these weakly dominated strategies are eliminated and all agents announce the true preferences about themselves, $R_i^i = \hat{R}_i$, it is strictly dominated to announce untruthful preferences about your neighbors, $R_{i+1}^i \neq \hat{R}_{i+1} = R_{i+1}^{i+1}$ or $R_{i-1}^i \neq \hat{R}_{i-1} = R_{i-1}^{i-1}$, since disagreeing with your neighbors is punished with the zero consumption bundle.

These two facts establish the main theorem in Sjöström (1994).

PROPOSITION 0. *Let f be an arbitrary social choice function. The mechanism described above implements f in UNE and in IWD.*

It is important to notice, for the discussion we undertake below, that the set of states in which not announcing the true preferences about oneself is weakly dominated are themselves states that typically produce very bad outcomes for other opponents (at least one of them will have zero consumption and probably many). If agents learn fast to avoid these (totally inconsistent) states, there is no incentive to tell the truth about oneself. The mechanism we have just described focuses on consensus announcements, since disagreement is punished with zero consumption; truth-telling is only rewarded in a set of states which need not be very prominent in the minds of the players. This is precisely the reason why, if agents are boundedly rational in the way we describe, convergence to the social choice outcome function may fail to occur.

We now move on to the characterization of the evolutionary dynamics we analyze.

Fix a given mechanism Γ and a given preference profile $R \in \Phi$. Let $x_i^{m_i}$ be the probability assigned by agent i to message m_i , and let $x_i \in \Delta_i$ be a mixed strategy for agent i (where Δ_i denotes the $|M_i| - 1$ -dimensional simplex which describes player i 's mixed strategy space). Let also $x_{-i} \in \Delta_{-i} \equiv \times_{j \neq i} \Delta_j$ be a mixed strategy profile for agents other than i , with $x \equiv (x_i, x_{-i}) \in \Delta \equiv \times_{i \in I} \Delta_i$. Finally, let $u_i(x_i, x_{-i}) = \sum_{m \in M} v_i(\alpha_i(m_i, m_{-i}), R_i) \Pi_{j \in I} x_j^{m_j}$.

We formalize player i 's behavior in terms of the mixed strategy $x_i(t)$ he or she adopts at each point in time. The vector $x(t)$ will then describe the *state of the system* at time t , defined over the state space Δ , with Δ^0 denoting its relative interior, i.e., the set of completely mixed strategy profiles.

Assumption d.1. The evolution of $x(t)$ is given by a system of continuous-time differential equations:

$$\dot{x}_i^{m_i} = D_i^{m_i}(x(t)). \quad (1)$$

We require that the autonomous system (1) satisfies the standard regularity condition; i.e., D must be (i) Lipschitz continuous with (ii) $\sum_{m_i \in M_i} D_i^{m_i}(x(t)) = 0$. Furthermore, D must also satisfy the following requirements:

Assumption d.2. D is a regular (payoff) *monotonic* selection dynamic. More explicitly, let $g_i(m_i, x_{-i}(t)) \equiv \dot{x}_i^{m_i}(t)/x_i^{m_i}(t)$ denote the growth rate of strategy m_i . Then for all $m_i, m'_i \in M_i$ and all $x_{-i} \in \Delta_{-i}$ it must be that

$$\begin{aligned} & \text{sign}[g_i(m_i, x_{-i}(t)) - g_i(m'_i, x_{-i}(t))] \\ &= \text{sign}[u_i(m_i, x_{-i}(t)) - u_i(m'_i, x_{-i}(t))]. \end{aligned}$$

Assumption d.2 is commonly used in the literature to capture the essence of a *selective* evolutionary process.⁹ Given the mixed strategy profile played at each point in time, strategies with higher expected payoff grow faster than poorly performing ones.

Assumption d.3. $x(0) \in \Delta^0$.

Assumption d.3 is also standard in the evolutionary literature. It excludes the possibility that the selection dynamic acts only on a subset of the strategy space. This possibility arises because any solution of a monotonic selection dynamics leaves (any face of) Δ , as well as Δ^0 , invariant (and, a fortiori, forward invariant). In other words, a strategy that has zero weight at time zero would also have zero weight at all subsequent times. If Assumption d.3 did not hold, the selection dynamics would then operate on a different game.

3. AN EXAMPLE

We prefix the dynamic analysis of the mechanism with the following example, taken from Sjöström (1994, p. 504), which is intended to convey the essence of our results. There is one unit of a single divisible private good, which has to be divided among three players: 1, 2, and 3. Preferences of players 1 and 2 are increasing in the amount of the good they consume, and are common knowledge for all players and the planner. There are two possible types for player 3's preferences, which are indexed by 0 and 1. Preferences of type 0 peak at consumption 1/3; preferences of type 1 peak at consumption 1/2. Player 3's type is common knowledge among the players, but the planner does not know it.

⁹ See, for example, Samuelson and Zhang (1992) and Weibull (1995).

m_2^0		m_2^1	
$\frac{1}{4}, \frac{1}{4}, \frac{1}{2}$	$\frac{1}{3}, 0, \frac{1}{3}$	$0, 0, \frac{1}{2}$	$0, \frac{1}{3}, \frac{1}{2}$
$0, \frac{1}{3}, \frac{1}{3}$	$0, 0, \frac{1}{3}$	$\frac{1}{3}, 0, \frac{1}{2}$	$\frac{1}{3}, \frac{1}{3}, \frac{1}{3}$
m_3^0		m_3^1	

FIG. 1. Sjöström’s example: game G .

For preferences of type 0, the social choice function recommends the outcome $f(0) = (1/4, 1/4, 1/2)$; for preferences of type 1, $f(1) = (1/3, 1/3, 1/3)$. Notice that the social choice function is such that type 3 would prefer the outcome $f(1)$ when she is of type 0, and the outcome $f(0)$ when she is of type 1. This provides her with an incentive to conceal her type, and therefore the planner needs a nontrivial mechanism to elicit her true preferences.

The mechanism proposed by Sjöström requires the three players to make a simultaneous statement about the preferences of player 3. Let $m_i^1(m_i^0)$, $i \in I$ represent the message in which preferences of type 1 (type 0) for player 3 are announced by player i . Figure 1 illustrates the outcome function of the mechanism. As for its dynamic analysis, we shall focus on the case in which true preferences of player 3 are of type 1, and assume that Fig. 1 also represents the game’s payoffs when player 3’s preferences are of type 1. We denote this game by G .

Player 1 picks a row, player 2 a column, and player 3 picks a matrix. We first note that the mechanism leads to a game which is *weakly dominance solvable*, as it can be reduced to a single outcome (the *solution*) by the iterated deletion of weakly dominated strategies. In particular, only one round of deletion of weakly dominated strategies, and then an additional round of deletion of strictly dominated strategies, is needed. In consequence, unlike other weakly dominance solvable games, the same outcome is selected independently on the order by which strategies are deleted.¹⁰ We start by deleting the weakly dominated strategy m_3^0 for player 3 (the other agents have no dominated strategies at this stage). The reason is that, like in the mechanism described in Section 2, truth-telling about your

¹⁰ Using Marx and Swinkels’ (1997) terminology, game G is in fact weakly dominance solvable in a *nice* sense. So is Sjöström’s (1994) general mechanism presented in the previous section.

own preferences never hurts, and is strictly optimal when the opponents disagree on your own type. Once m_3^0 has been removed, strategies m_1^0 and m_2^0 become strictly dominated. The reason is that, like in the mechanism described in Section 2, if all the players tell the truth about their own preference, lying about a neighbor is punished with zero consumption. The unique strategy profile selected is then (m_1^1, m_2^1, m_3^1) , that is, the pure strategy profile in which the true preferences are consistently revealed.

Since each player has only two pure strategies in her support, we abuse our notation setting $x_i \equiv x_i^{m_i^1}$.¹¹ In the following proposition we fully characterize the set of Nash equilibria of game G .

PROPOSITION 1. *The set NE of Nash equilibria of G is the union of precisely two disjoint components NE^0 and NE^1 , where*

$$NE^0 \equiv \{x \in \Delta \mid x_1 = x_2 = 0, x_3 \leq \frac{3}{7}\},$$

$$NE^1 \equiv \{x \in \Delta \mid x_1 = x_2 = 1, x_3 \geq \frac{1}{2}\}.$$

Proof. See Cabrales and Ponti (1998, Proposition 1). ■

We now move on to dynamics. Denote by $RE(G)$ the set of restpoints of G under any monotonic dynamic. It is straightforward to show that $RE(G)$ contains (together with all the pure strategy profiles) only the components

$$RE^0 \equiv \{x \in \Delta \mid x_1 = x_2 = 0, x_3 \in [0, 1]\}$$

and

$$RE^1 \equiv \{x \in \Delta \mid x_1 = x_2 = 1, x_3 \in [0, 1]\}.$$

Our task is to study the asymptotics of a monotonic selection dynamic whose initial state lies in the relative interior of the state space.

PROPOSITION 2. *Any interior solution $x(t, x(0))$ of a monotonic selection dynamics $\dot{x} = D(x)$ converges to NE.*

Proof. See the Appendix. ■

If initial conditions are completely mixed, we then know from Proposition 2 that the evolutionary dynamics will eventually converge to a Nash equilibrium of the game. In the next section we show that this result generalizes locally also in the case of Sjöström's (1994) mechanism, as described in Section 2.

¹¹ The fact that each player has only two available options also allows us to express the dynamics in terms of the payoff difference between player i 's truthful and untruthful strategy, which we call $\Delta\Pi_i(x_{-i}(t))$ (i.e., $\Delta\Pi_i(x_{-i}(t)) \equiv u_i(m_i^1, x_{-i}(t)) - u_i(m_i^0, x_{-i}(t))$).

4. LOCAL RESULTS FOR THE GENERAL GAME

In Proposition 3 we characterize some components of Nash equilibria for the game induced by the mechanism. In particular, we show that any message profile in which the agents are unanimous in the (arbitrary) preference profile they announce, R^* (more precisely, the preferences they announce about their neighbors and themselves are taken from R^*), is an equilibrium. Furthermore, any mixed strategy profile in which agents mix between messages consistent with R^* and other messages that only differ in the announcements agents make about their own preferences is also an equilibrium, provided that messages in R^* are given a high enough weight. The equilibria in each of these components are not payoff equivalent, since disagreeing with a neighbor (event with nonzero probability in these mixed equilibria) results in a punishment. Nevertheless, Proposition 4 shows that this punishment is not high enough to prevent these equilibria to be the limit points of some interior path of any monotonic selection dynamic.

Before we proceed, some further terminology is needed. Let $m_i^* = (R_{i-1}^*, R_i^*, R_{i+1}^*)$ be a consensus announcement by agent i , let $U_i = \max_R v_i(f_i(R), \hat{R}_i)$ be the utility associated to the most preferred outcome from the social choice function for agent i with true preferences \hat{R}_i , and let $U_{in} = \max_R v_i(\frac{1}{n}f_i(R), \hat{R}_i)$ be the utility associated to the most preferred consumption bundle among those that result from dividing the bundles assigned by the social choice function by n . Let also $S_i(R^*)$ denote the set of all pure strategies in which announcements about the neighbors agree with R^* , i.e.,

$$S_i(R^*) = \{m_i \in M_i \mid R_{i-1}^i = R_{i-1}^*, R_{i+1}^i = R_{i+1}^*\}, \quad (2)$$

with $\bar{S}_i(R^*) = \{m_i \in M_i \mid m_i \notin S_i(R^*)\}$ denoting the complement of $S_i(R^*)$ with respect to M_i and $S_{-i}(R^*) \equiv \times_{j \neq i} S_j(R^*)$ ($\bar{S}_{-i}(R^*) \equiv \times_{j \neq i} \bar{S}_j(R^*)$). Finally, denote by $S_i^{k_i}(R^*)$ the following

$$S_i^{k_i}(R^*) = \{x_i \mid x_i^{m_i} = 0, \text{ for all } m_i \notin S_i(R^*) \text{ and } x_i^{m_i^*} > k_i\}, \quad (3)$$

where we assume

$$(k_i)^n \geq \frac{U_{in} - v_i(0, \hat{R}_i)}{v_i(f_i(\phi(i, R^*)), \hat{R}_i) - v_i(0, \hat{R}_i) + U_{in} - v_i(0, \hat{R}_i)} \quad (4)$$

for all i and $j \neq i$. The set $S_i^{k_i}(R^*)$ is the set of all mixed strategies in which i 's announcements about her neighbors agrees with R^* , and the probability of announcing R_i^* is higher than k_i .

PROPOSITION 3. For all $\hat{R}, R^* \in \mathfrak{R}$ and $x_i \in S_i^{k_i}(R^*)$, x is a Nash equilibrium of (Γ, \hat{R}) .

Proof. See the Appendix. ■

To understand the role of (4) in the proof of Proposition 3, notice that, against any $x_{-i} \in \times_{j \neq i} S_j^{k_j}(R^*)$, the payoff for agent i using strategy $m_i \in S_i(R^*)$ satisfies the following condition:

$$\begin{aligned} u(m_i, x_{-i}) &\geq \left(\min_{j \neq i} k_j \right)^{n-1} v_i(f_i(\phi(i, R^*)), \hat{R}_i) \\ &\quad + \left(1 - \left(\min_{j \neq i} k_j \right)^{n-1} \right) v_i(0, \hat{R}_i). \end{aligned} \quad (5)$$

This is because, for all $j \neq i$, $x_j^{m_j} \geq k_j$, which in turn implies a lower bound (i.e., $(\min_{j \neq i} k_j)^{n-1}$) on the probability with which m_{-i} is *totally consistent* with $m_i \in S_i(R^*)$ and, therefore, the payoff $v_i(f_i(\phi(i, R^*)), \hat{R}_i)$ is achieved. With the remaining probability $1 - (\min_{j \neq i} k_j)^{n-1}$, the worse that can happen to player i is that her message does not match the announcements of his or her neighbors about themselves, in which case his or her payoff is $v_i(0, \hat{R}_i)$. By the same token, against any $x_{-i} \in \times_{j \neq i} S_j^{k_j}(R^*)$, the payoff for agent i announcing a message $m'_i \in \bar{S}_i(R^*)$ is at most

$$u(m'_i, x_{-i}) \leq \left(\min_{j \neq i} k_j \right)^{n-1} v_i(0, \hat{R}_i) + \left(1 - \left(\min_{j \neq i} k_j \right)^{n-1} \right) U_{in}. \quad (6)$$

From Eqs. (5) and (6), it follows that

$$\begin{aligned} &u_i(m_i, x_{-i}) - u_i(m'_i, x_{-i}) \\ &\geq v_i(0, \hat{R}_i) - U_{in} + \left(\min_{j \neq i} k_j \right)^{n-1} \\ &\quad \times (v_i(f_i(\phi(i, R^*)), \hat{R}_i) + U_{in} - 2v_i(0, \hat{R}_i)), \end{aligned} \quad (7)$$

which implies $u_i(m_i, x_{-i}) - u_i(m'_i, x_{-i}) \geq 0$, provided that (4) is satisfied.

Also note that, for all $x_{-i} \in \times_{j \neq i} S_j^{k_j}(R^*)$, if $m_i, m'_i \in S_i(R^*)$, then $u_i(m_i, x_{-i}) - u_i(m'_i, x_{-i}) = 0$. This is because, in playing any strategy in $S_i(R^*)$, agent i rules out the possibility that totally inconsistent states occur (at least the announcements about i have to coincide). These are the only states in which i 's announcement about her own preferences makes a difference to her own payoff.

We shall now prove the elements in all the Nash equilibria components characterized by Proposition 3 are *reachable*, i.e., are limit points for some interior solution. By Lipschitz continuity, there exists a constant $K > 0$ such that for all m_i , x_{-i} , and x'_{-i} , we have that

$$|g_i(m_i, x_{-i}(t)) - g_i(m_i, x'_{-i}(t))| \leq K|x_{-i} - x'_{-i}|,$$

where the $|\cdot|$ denotes the norm of a vector. This in turn implies that, for all $h_v > 0$ with $u_i(m_i, x_{-i}(t)) - u_i(m'_i, x_{-i}(t)) \leq -h_v$, there exists some $h_g > 0$, such that $g_i(m_i, x_{-i}(t)) - g_i(m'_i, x_{-i}(t)) \leq -h_g$. By analogy with (7), for any $m_i \in \bar{S}_i(R^*)$, it also must be

$$\begin{aligned} & u_i(m_i, x_{-i}) - u_i(m_i^*, x_{-i}) \\ & < U_i - v_i(0, \hat{R}_i) \\ & \quad - \Pi_{j \neq i} x_j^{m_j^*}(t) \left(v_i(f_i(\phi(i, R^*)), \hat{R}_i) + U_i - 2v_j(0, \hat{R}_i) \right). \end{aligned}$$

Therefore, if h_v is a constant such that $0 \leq h_v < \min_{i,R} v_i(f_i(\phi(i, R^*)), \hat{R}_i) - v_i(0, \hat{R}_i)$, then there exists another constant $H \in [0, 1)$, with

$$H = \max_i \left\{ \left(\frac{U_i - v_i(0, \hat{R}_i) + h_v}{v_i(f_i(\phi(i, R^*)), \hat{R}_i) + U_i - 2v_i(0, \hat{R}_i)} \right)^{1/(n-1)} \right\},$$

such that, if $x_j^{m_j^*}(t) > H$ for all j and t , then strategies not in $S_i(R^*)$ are decreasing at a rate not higher than $-h_g$.

We also need to establish a link between the weight with which messages $m_{-i} \in \bar{S}_{-i}(R^*)$ are played and the relative performance of strategies $m_i \in S_i(R^*)$. This is done by means of the function

$$X_i(t) = \sum_{j \neq i} \left(\left(\sum_{m_j \in \bar{S}_j} x_j^{m_j}(t) \right)^2 + \sum_{m_j \in \bar{S}_j} (x_j^{m_j}(t))^2 \right),$$

with $X(t) = \max_i [X_i(t)]$. The function $X_i(t)$ accounts for the relative weight of messages $m_{-i} \in \bar{S}_{-i}(R^*)$ in x_{-i} , since only against these messages do strategies in $S_i(R^*)$ yield different payoffs for player i . Therefore, the maximum difference in payoffs between strategies in $S_i(R^*)$, and

therefore in growth rates by monotonicity, is connected to $X_i(t)$, as shown in Lemma 2.¹² Finally, let

$$L = \min_i \left\{ \exp \left(\frac{-KX(0)}{h_g} \frac{H^2}{(x_i^{m_i^*}(0))^2} \right) \right\}.$$

The constant L appears because we want to show that $x_i^{m_i^*}(t)$ need not go to one in the limit, even if there is convergence to the equilibrium component to which m^* belongs. For any $m_i \in S_i(R^*)$, the ratio

$$\frac{x_i^{m_i^*}(t)}{x_i^{m_i}(t)} \frac{x_i^{m_i}(0)}{x_i^{m_i^*}(0)}$$

is the integral of the differences in growth rates (thus connected to the difference in payoffs by monotonicity) between m_i and m_i^* . This integral depends on $X(t)$, as we show in Lemma 2. But $X(t)$ depends on $X(0)$ also, as well as on the growth rates of strategies of i 's opponents in $\bar{S}_{-i}(R^*)$. As shown in the following Proposition 4, also the weight of these latter strategies has an upper bound which depends on h_g , K , and H . Thus, the constant L can be used to set an upper bound for the integral of the difference in growth rates between strategies m_i and m_i^* .

Also notice that $X(0)$ can be made arbitrarily close to zero (and, therefore, L arbitrarily close to 1) by selecting an initial condition in which the aggregate weight of strategies in $\bar{S}_{-i}(R^*)$ is arbitrarily small.

PROPOSITION 4. *Assume that, for all $i \in I$, $x_i^{m_i^*}(0)L > H$. Then*

(a) *for all $m_i \in \bar{S}_i(R^*)$, $x_i^{m_i}(t)/x_i^{m_i}(0) < \exp[-h_g t]H/x_i^{m_i^*}(0)$ for all t and all i ;*

(b) *$x_i^{m_i^*}(t) > H$ for all t ;*

(c) *$x_i^{m_i^*}(t)/x_i^{m_i}(t) < (x_i^{m_i^*}(0)/x_i^{m_i}(0))(1/L)$ for all t and $m_i \in S_i(R^*)$.*

Proof. See the Appendix. ■

Proposition 4(b) guarantees that pure strategy equilibria in all equilibrium components (including the “wrong” ones) are attractors of interior paths. By Proposition 4(c), the limiting weight of m_i^* is less than 1 (provided L is sufficiently close to 1), and therefore some mixed strategy equilibria are attractors as well, if the initial conditions give sufficiently little weight to strategies in $\bar{S}_{-i}(R^*)$. This guarantees that, even if there is convergence to the “right” component, it need not be to the pure strategy

¹² In the Appendix.

equilibrium (remember that the equilibria are not payoff equivalent, as the mixed strategy equilibria have lower expected payoff because agents are punished for announcing discordant preferences).

Convergence to mixed equilibria may occur because payoffs to all strategies in $S_i(R^*)$ are “close,” if the weight of strategies in $\bar{S}_{-i}(R^*)$ is small. We know, by Proposition 4(a), that the weight of strategies in $\bar{S}_{-i}(R^*)$ is indeed decreasing. So, even though m_i^* has a payoff advantage, this advantage vanishes, and Assumption d.2 (plus Lipschitz continuity) guarantees that it does not accumulate fast enough.

5. MORE ON THE EXAMPLE (STABILITY WITH/OUT DRIFT)

In the previous section, we extended the convergence result of Proposition 2 to the general mechanism, showing that the limit points of the dynamics for interior initial conditions are generally different from the outcomes intended by the planner. We now go back to our example to test the stability properties of the various equilibrium components.

DEFINITION 1. Let $x(t, x(0))$ be the solution of (1) on state space Δ given initial conditions $x(0)$. Let also C be a closed set of restpoints in Δ of the same differential equation. Then:

(i) C is (interior) *stable* if, for every neighborhood O of C , there is another neighborhood U of C , with $U \subset O$, such that for any $x(0) \in U \cap \Delta(U \cap \Delta^0)$ we have $x(t, x(0)) \in O$;

(ii) C is (interior) *attracting* if it is contained in an open set O such that for any $x(0) \in O \cap \Delta(O \cap \Delta^0)$ we have $\lim_{t \rightarrow \infty} x(t, x(0)) \in O$;

(iii) C is *globally* (interior) *attracting* if for any $x(0) \in \Delta(\Delta^0)$ we have $\lim_{t \rightarrow \infty} x(t, x(0)) \in O$;

(iv) C is called (interior) *asymptotically stable* if it is (interior) attracting and (interior) stable.

To simplify the analysis, we set additional conditions on the dynamics, which is the purpose of the following assumption, (which replaces Assumptions d.1–3):

Assumption d.4. The evolution of $x(t)$ is given by the system of continuous-time differential equations

$$\dot{x}_i \equiv \tilde{D}_i(x(t), \lambda) = x_i(t)(1 - x_i(t)) \Delta \Pi_i(\cdot) + \lambda(\beta_i - x_i(t)), \quad (8)$$

with $\lambda \geq 0$, $\beta_1 = \beta_2 = \frac{1}{2}$, and $\beta_3 = \beta \in (0, 1)$.

In words, the evolutionary dynamic is now composed of two additive terms. The first represents the standard replicator dynamic, while the second term ensures that, at each point in time, each strategy is played with positive probability, no matter how it performs against the current opponents' mixed strategy profile (i.e., it points the dynamic *toward the relative interior* of the state space Δ). Following Binmore and Samuelson (1999), this latter term is called *drift*: it opens the model to the possibility of a heterogeneity of behaviors. Gale et al. (1995) derive an analogous system in the following way. At each point in time, a fixed proportion of players (of measure $\frac{\lambda}{1+\lambda}$) is replaced by new individuals whose aggregate behavior is represented by a generic, constant, completely mixed strategy (i.e., β_i), while the rest of the population aggregate behavior follows the replicator dynamics. The relative importance of the drift is measured by λ , which we refer to as the *drift level*. We assume λ to be "very small," reflecting the fact that all the major forces which govern the dynamics should be captured by the evolutionary dynamic defined by D , which here takes the form of the replicator dynamics.

We check how the model reacts to the introduction of such a perturbation. The stability analysis of the replicator dynamics with drift will give us information about the effects of small changes in the vector field on the equilibria of the system defined by the replicator dynamic (in other words, it will test the *structural stability* of such equilibria). To simplify the exposition, β_1 and β_2 have been chosen to be $1/2$, since only the value of β_3 turns out to be genuinely significant.

We start by looking at the case of the replicator dynamic without drift (i.e., when $\lambda = 0$). We know from Proposition 2 that NE is globally interior attracting, since it attracts every interior path under any monotonic selection dynamic (of which the replicator dynamic is a special case). We now take a closer look at the stability properties of each component of Nash equilibria separately (i.e., NE^0 and NE^1).

Figure 2 shows a phase diagram describing trajectories of the replicator dynamic starting from some interior initial conditions. The Nash equilibrium component NE^0 (NE^1) is represented by a bold segment in the bottom-left (top-right) corner of the state space Δ . First notice that, as we know from Proposition 2, all trajectories converge to a Nash equilibrium of the game. Moreover, the diagram shows (consistently with Proposition 4) that there are some trajectories of the replicator dynamic which converge to NE^0 , the Nash equilibrium component in which both players 1 and 2 deliver the false message with probability 1. However, this latter component is not asymptotically stable, as can be easily spotted from the diagram. Trajectories starting arbitrarily close to NE^0 , provided $x_3 > \frac{3}{7}$, will eventually converge to the truth-telling component. We summarize the

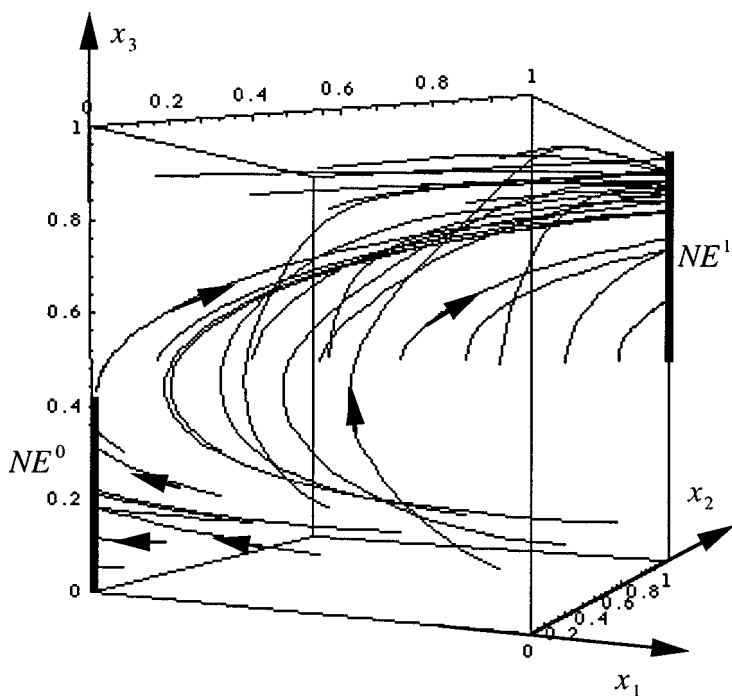


FIG. 2. The replicator dynamic and game G .

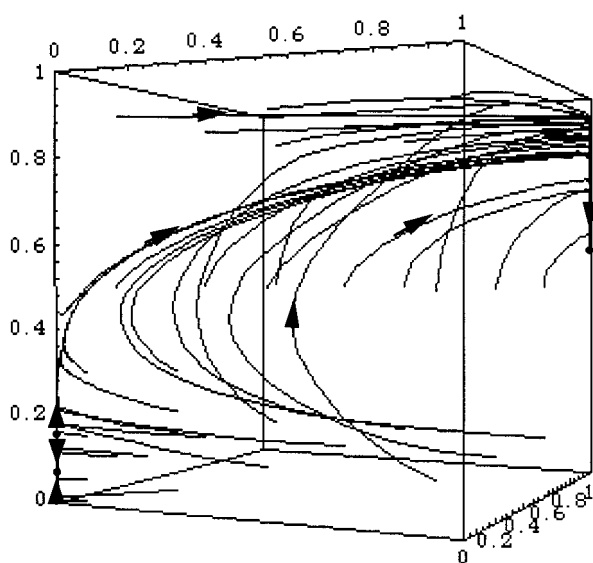
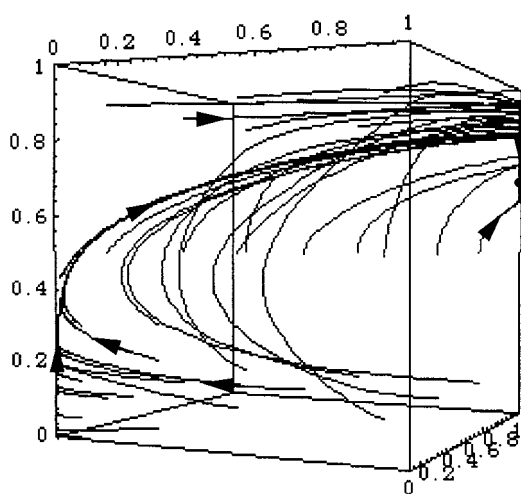
key properties of these trajectories in the following proposition:

PROPOSITION 5. *Under the replicator dynamic (i) NE^1 is interior asymptotically stable, whereas (ii) NE^0 is not.*

Proof. See Cabrales and Ponti (1998, Proposition 5). ■

We now move to the analysis of the replicator dynamic with drift.

Let $\beta \in (0, 1)$ be a generic element of the space of the feasible perturbation. Figure 3 shows trajectories of the replicator dynamic with drift under two different specifications of β . Figure 3(b) represents a situation in which, in the proximity of NE^0 , the drift against m_i^0 is uniform across players, where in Fig. 3(a) the drift against m_3^0 is lower. As the diagrams show, there is a local attractor close to NE^1 in both cases. Moreover, none of the elements of NE^0 is a restpoint of the dynamic with drift in Fig. 3(b). In contrast, in Fig. 3(a) there is an additional local attractor which belongs to NE^0 : trajectories starting close to NE^0 converge to it, as it happens in the case of the replicator dynamics without drift.

a) $\beta=1/100$ b) $\beta=1/2$ FIG. 3. The dynamic with drift and game G .

We are interested in the convergence and stability properties of (8) when $\lambda \rightarrow 0$, considering two different configurations of the drift parameter β :

$$\text{CASE A:} \quad \beta \in \left(0, \frac{23 - 4\sqrt{30}}{49}\right)$$

$$\text{CASE B:} \quad \beta \in \left(\frac{23 - 4\sqrt{30}}{49}, 1\right).$$

Given $\frac{23 - 4\sqrt{30}}{49} \approx 0.0222673$, CASE A depicts a situation in which, for small values of x_i , the drift against the untruth-telling strategy is substantially lower for player 3 than for her opponents.

In the following proposition we characterize the set of restpoints of the dynamic with drift, together with their stability properties:

PROPOSITION 6. *Let $\hat{\text{RE}}(\beta)$ be the set of restpoints of (8) for λ sufficiently close to 0. The following properties hold:*

(a) $\forall \beta \in (0, 1)$, $\hat{\text{RE}}(\beta)$ contains an element of NE^1 , which is also asymptotically stable.

(b) Under CASE A $\hat{\text{RE}}(\beta)$ contains also two additional restpoints, both belonging to NE^0 , one of which is asymptotically stable.

Proof. See the Appendix. ■

As we acknowledged in Section 1, there is a striking similarity between the content of Proposition 6 and the findings of Gale et al. (1995). They also find that, for the *entry game*, the Nash component in which the incumbent carries out her “incredible threat” is reachable under the replicator dynamics. Moreover, like our NE^0 , it fails to be interior asymptotically stable, but for certain parameter values it may be asymptotically stable when the system is slightly perturbed. Given the failure of asymptotic stability without perturbations, one would expect any perturbation to move the system away from the unstable component and the weakly dominated strategy to become extinct. Proposition 6 tells us that evolutionary game theory does not provide a ground for such a claim. Once again, the intuition here is similar to the one in Gale et al. (1995). When there is drift, strategies against which the weakly dominated strategy does poorly will have positive weight at all times and, therefore, the part of the dynamics that depends on payoffs pushes against the dominated strategy. On the other hand, drift may provide a direct push in favor of the dominated strategy (and more crucially, in favor of those strategies of the other players which do well against the dominated strategy). When the

balance between these two forces is right, one gets a stable equilibrium with non-negligible weight for the dominated strategy.

6. BEST-REPLY DYNAMICS AND SJÖSTRÖM'S MECHANISM

In this section, we consider an alternative scenario. Suppose that $x(t)$ evolves according to the dynamics

$$\dot{x} = \text{BR}(x) - x, \quad (9)$$

with $\text{BR}(x)$ denoting the *mixed strategy best-reply correspondence* $\text{BR}: \Delta \mapsto \Delta$.¹³ This alternative dynamic defines a (continuous-time) version of the classic *best-reply* dynamics, often proposed as an alternative to the evolutionary dynamics studied hereto. We can give two interpretations to (9). Following Matsui (1992), we can use (9) to approximate the evolution of an infinite population of players who occasionally update their strategy, selecting a best reply to the current population state $x(t)$. Alternatively, (9) can be regarded as the continuous-time limit (up to a reparametrization of time) of the well known *fictitious play* dynamic.¹⁴ This dynamic accounts for the evolution of players' *beliefs*, when these beliefs follow the empirical frequencies with which each pure strategy profile has been played (and perfectly observed) in the past, and agents select, at each point in time, a pure strategy among those which maximize their expected payoff, given their current beliefs.

Let $\hat{S}_i = \{m_i \in M_i \mid R_i^i = \hat{R}_i^i\}$, $\hat{s}_i = \{m_i \in \hat{S}_i \mid R_{i-1}^i = \hat{R}_{i-1}^i, R_{i+1}^i = \hat{R}_{i+1}^i\}$, with $\hat{s} = \{\hat{s}_i\}$, $i \in I$ denoting the pure Nash equilibrium in which all agents consistently reveal their true preferences (i.e., the "solution" of Γ given the true preference profile \hat{R}).

PROPOSITION 7. *Any interior solution of (9) converges to \hat{s} .*

Proof. See the Appendix. ■

¹³ Notice that, for some $x \in \Delta$, $\text{BR}(x)$ can take infinitely many values. Thus, uniqueness of the solution of (9) is not guaranteed. However, since $\text{BR}(x)$ is upper-hemicontinuous with closed and convex values, it can be shown that the differential inclusion $\dot{x} \in \text{BR}(x) - x$ has at least one (interior) solution $x(t, x(0))$, which is Lipschitz continuous and defined, for any $t \geq 0$ (Aubin and Cellina, 1984, Chap. 2). On the stability properties of (9) see Hofbauer (1997).

¹⁴ First introduced by Brown (1951) as an algorithm to compute Nash equilibria, fictitious play has been recently re-interpreted as a learning model by Fudenberg and Kreps (1993). We prefer here the non-standard version in continuous-time to be consistent with the rest of the paper. Nevertheless, in an earlier version of this paper we prove that the same results still hold if the dynamics are defined in discrete-time.

For best-reply dynamics we have shown that every interior solution converges to the unique equilibrium whose outcome is the one the planner wants to implement. This is so because since initial beliefs are completely mixed, they will always be completely mixed, so these weakly dominated strategies will always remain suboptimal, will never be played, and their weight in beliefs will eventually vanish. This implies that nonequilibrium strategies by which agents misrepresent their neighbors' preferences become also suboptimal, and agents will learn not to use them.¹⁵

The results obtained here are so different from those we derived in the previous sections essentially because the difference in growth rates between two pure strategies, in the case of the best-reply dynamics (9), need not satisfy Lipschitz continuity. The only strategies with a positive growth rate are best responses; this implies that there is an infinite response of growth rates to changes in the sign of the differences in payoffs, which is precisely what Lipschitz continuity rules out.

7. CONCLUSIONS

We have argued that there is room for doubt about the practicability of one of the leading examples of implementation with iterated deletion of weakly dominated strategies when agents are boundedly rational. As we said in Section 1, there are only few papers that study implementation with boundedly rational players, so a deeper theoretical study with evolutionary tools of other mechanisms studied in the literature would enhance our understanding of the performance of these mechanisms with this type of agent, a necessary step before mechanisms are used in real life.

Further empirical study is at least as necessary. It would, for example, help to answer the question about which of the dynamics assumptions is more appropriate. In this sense, there is already some evidence on mechanism design and learning algorithms. Chen and Tang (1998) have done experiments with the basic quadratic mechanism by Groves and Ledyard (1977) and the paired-difference mechanism by Walker (1981). They estimate different learning models using experimental data, showing that variants of stimulus-response learning algorithms (whose expected law of motion is the replicator dynamics) outperform the generalized fictitious play model. This is also consistent with the good performance that Roth and Erev (1995) show for stimulus-response learning algorithms in mim-

¹⁵ By analogy, it can be proved that every interior trajectory of (9), for game Γ , converges to $(1, 1, 1)$ (see Cabrales and Ponti, 1998, Proposition 7).

icking the behavior of a range of experimental data, which includes other weakly dominance solvable games, like the ultimatum game.¹⁶

But even more importantly, the empirical and experimental work would help to design games with good convergence properties to the preferred social outcome by revealing how people adjust their play in games like that studied in this paper, as well as in other mechanisms proposed by the literature. We have already begun to do such experimental studies.¹⁷

APPENDIX

Proof of Proposition 2. To prove the proposition, it is enough to show that any interior trajectory converges. The reason is that, once convergence has been proved, we can apply the standard result “convergence implies Nash under any monotonic selection dynamics” (see, e.g., Weibull, 1995, Theorem 5.2 (iii)).

We start by observing that the dynamic is forward invariant. This implies that $x_i(t)$ is always defined and positive, for any nonnegative t . By monotonicity, $x_3(t)$ is also a positive, increasing function of t and bounded above by 1 (since m_3^1 is a weakly dominant strategy). Therefore, $x_3(t)$ must converge (this already implies convergence of player 3’s mixed strategy). Let $x_i^* \equiv \lim_{t \rightarrow \infty} x_i(t)$, when such a limit exists. Three alternative cases have to be discussed.

(a) $x_3^* = 0$. If $x_3^* = 0$ there must be a time t' such that $x_3(t) < \frac{3}{7}$ for $t > t'$. This implies that there is a $k > 0$ such that for all $t' > t$, $\Delta \Pi_i(x(t)) < -k$ for $i = 1, 2$. This implies, by monotonicity, $\lim_{t \rightarrow \infty} x_i(t) = 0$ for $i = 1, 2$; thus $x^* = (0, 0, 0)$.

(b) $x_3^* = 1$. By a similar argument, monotonicity implies $x^* = (1, 1, 1)$.

(c) $x_3^* \in (0, 1)$. We want to prove that x_3^* cannot converge to a value within this range unless the system converges to a Nash equilibrium. To do so (given the special features of our example) it is enough to show that, if $x_3^* \in (0, 1)$, then both players 1 and 2 select, in the limit, the same pure strategy. Given that this implies convergence of the full mixed strategy

¹⁶ In their paper, Roth and Erev (1995) show that these dynamics explain the data significantly better, according to quadratic deviation measures and others, than a generalized fictitious play model which can accommodate behaviors ranging from fictitious play to best response dynamics by the estimation of a “forgetfulness parameter” which weights past information. For the experimental evidence on learning rules, see also Tang (1996), Chen et al. (1997), and Mookherjee and Sopher (1997).

¹⁷ See Cabrales et al. (1998) and Gantner et al. (1998).

profile, the result follows. More formally, what we need to prove is contained in the following lemma:

LEMMA 1. *If $x_3^* \in (0, 1)$ then either*

$$x_i^* = 0, \quad i = 1, 2 \quad (\text{CASE 0 hereafter})$$

or

$$x_i^* = 1 \quad i = 1, 2 \quad (\text{CASE 1}).$$

Proof. Assume, for the purpose of contradiction, that neither of the above statements is true. In this case, there must exist a sequence $\{t_k\}_{k=1}^\infty$ and a positive constant $\varepsilon > 0$ such that either $x_i(t_k) > \varepsilon$, $i = 1, 2$, or $x_i(t_k) < 1 - \varepsilon$, $i = 1, 2$ for all k (in other words, assume that the system stays infinitely often an ε away from the faces of Δ in which player 1 and 2 play the same pure strategy). We already noticed that these are the only faces of Δ in which both pure strategies for player 3 yield the same payoff. If the system stays away from these faces infinitely often along the solution path, then the integral of the payoff difference $\Delta\Pi_3(x(t))$ goes to infinity as t goes to infinity.

To show this, notice that $\Delta\Pi_i(x(t))$ is a continuous function of $x(t)$ defined over a compact set (Δ). In the case of player 3, such a function takes the following form:

$$\Delta\Pi_3(x(t)) \equiv \frac{(x_1(t) - x_2(t))^2 + x_1(t)(1 - x_1(t)) + x_2(t)(1 - x_2(t))}{6}. \quad (10)$$

Take $g_M \equiv \max_{i \in I, x_{-i} \in \Delta_{-i}} [g_i(m_i, x_{-i}(t))]$, i.e., the highest possible growth rate (in absolute value) over all strategies and players (we know a max exists, since also $g_i(\cdot)$ is continuous in Δ). Then define τ_1 , τ_2 , τ_3 , and τ_4 as follows:

$$\begin{aligned} \tau_1 \quad \text{solves} \quad \varepsilon \exp[-g_M \tau_1] &= \frac{\varepsilon}{2} & \left(\text{i.e., } \tau_1 = \frac{\ln[2]}{g_M} \right), \\ \tau_2 \quad \text{solves} \quad (1 - \varepsilon) \exp[-g_M \tau_2] &= \frac{\varepsilon}{2} & \left(\text{i.e., } \tau_2 = \frac{\ln\left[-2 + \frac{2}{\varepsilon}\right]}{g_M} \right), \end{aligned}$$

$$\begin{aligned} \tau_3 \quad \text{solves} \quad \varepsilon \exp[g_M \tau_3] &= 1 - \frac{\varepsilon}{2} & \left(\text{i.e., } \tau_3 = \frac{\ln \left[-\frac{1}{2} + \frac{1}{\varepsilon} \right]}{g_M} \right), \\ \tau_4 \quad \text{solves} \quad (1 - \varepsilon) \exp[g_M \tau_4] &= 1 - \frac{\varepsilon}{2} & \left(\text{i.e., } \tau_4 = \frac{\ln \left[\frac{2 - \varepsilon}{2 - 2\varepsilon} \right]}{g_M} \right). \end{aligned}$$

Let $\partial\tau \equiv \min[\tau_1, \tau_2, \tau_3, \tau_4]$ be the lower bound for the time interval in which, after each t_k , $\frac{\varepsilon}{2} < x_i < 1 - \frac{\varepsilon}{2}$, $i = 1, 2$ and therefore $\Delta\Pi_3(x(t))$ still remains bounded away from 0 (i.e., $\Delta\Pi_3(x(t)) > \frac{\varepsilon(1-\varepsilon/2)}{3} > 0$, $\forall t \in [t_k, t_k + \partial\tau]$). Denote by $G_\varepsilon = \{x \in \Delta \mid \Delta\Pi_3(x) \geq \frac{\varepsilon(1-\varepsilon/2)}{3}\}$. Now define $\gamma_i(x(t))$ as the time derivative of the log of the ratio between the probabilities with which each of player i 's pure strategies are played, which can be expressed in terms of the difference in the growth rates:

$$\begin{aligned} \gamma_i(x(t)) &\equiv \frac{\partial}{\partial t} \ln \left(\frac{x_i(t)}{1 - x_i(t)} \right) = \frac{\dot{x}_i(t)}{x_i(t)} - \frac{(1 - \dot{x}_i(t))}{1 - x_i(t)} \\ &= \frac{\dot{x}_i(t)}{x_i(t) - (x_i(t))^2}. \end{aligned}$$

Also $\gamma_3(x(t))$ is a positive number bounded away from zero infinitely often since, by Assumption d.1, it is a continuous function of $x(t)$ defined on a compact set, which preserves the same sign of $\Delta\Pi_3(x(t))$. This implies that we can always define a constant $g_\varepsilon = \min_{x \in G_\varepsilon} \gamma_3(x(t))$, with $g_\varepsilon > 0$ by Assumption d.2. Also by Assumption d.2, $\gamma_3(x(t)) > g_\varepsilon \Leftrightarrow \Delta\Pi_3(x(t)) > \frac{\varepsilon(1-\varepsilon/2)}{3}$. If we integrate the value of $\gamma_3(x(t))$ over time we then obtain

$$\lim_{t \rightarrow \infty} \int_0^t \gamma_2(x(t)) dt \geq \sum_{k=1}^{\infty} \int_{t_k}^{t_k + \partial\tau} \gamma_3(x(t)) dt > g_\varepsilon \sum_{k=1}^{\infty} \int_{t_k}^{t_k + \partial\tau} dt = \infty,$$

which implies that $x_3^* = 1$, which leads to a contradiction. ■

To summarize, Lemma 1 shows that, if $x_3^* \in (0, 1)$, $x_1(t)$, and $x_2(t)$ must converge (and therefore $x(t)$ must converge to a Nash equilibrium). Since this exhausts all cases the result follows. ■

Proof of Proposition 3. We begin by noting that, against any $m_{-i} \in M_{-i}$, all strategies $m_i \in S_i(R^*)$ yield the same payoff, as they only differ in i 's announcement about herself. Since $\text{supp}[x_{-i}] \subseteq S_{-i}(R^*)$, totally inconsis-

tent states (the only states where announcements about i 's own type influence her own payoff) are excluded.

For all \hat{x}_i such that $\hat{x}_i^{m_i} > 0$ only if $m_i \in S_i(R^*)$ we have

$$u_i(\hat{x}_i, x_{-i}) \geq \Pi_{j \neq i} x_j^{m_j^*} v_i(f_i(\phi(i, R^*)), \hat{R}_i) + (1 - \Pi_{j \neq i} x_j^{m_j^*}) v_i(0, \hat{R}_i).$$

For all $\bar{x}_i \neq \hat{x}_i$,

$$\begin{aligned} u_i(\bar{x}_i, x_{-i}) &\leq \sum_{m_i \in S_i} \bar{x}_i^{m_i} u_i(\hat{x}_i, x_{-i}) + \left(1 - \sum_{m_i \in S_i} \bar{x}_i^{m_i}\right) \\ &\quad \times \left[\Pi_{j \neq i} x_j^{m_j^*} v_i(0, \hat{R}_i) + (1 - \Pi_{j \neq i} x_j^{m_j^*}) U_{in}\right]. \end{aligned}$$

Then

$$\begin{aligned} &u_i(\hat{x}_i, x_{-i}) - u_i(\bar{x}_i, x_{-i}) \\ &\geq \left(1 - \sum_{m_i \in S_i(R^*)} \bar{x}_i^{m_i}\right) \left[\Pi_{j \neq i} x_j^{m_j^*} (v_i(f_i(\phi(i, R^*)), \hat{R}_i) - v_i(0, \hat{R}_i)) \right. \\ &\quad \left. + (1 - \Pi_{j \neq i} x_j^{m_j^*}) (v_i(0, \hat{R}_i) - U_{in})\right], \end{aligned}$$

which is great than zero since, by (4),

$$\begin{aligned} \Pi_{j \neq i} x_j^{m_j^*} &\geq \Pi_{j \neq i} k_j \\ &\geq \frac{U_{in} - v_i(0, \hat{R}_i)}{v_i(f_i(\phi(i, R^*)), \hat{R}_i) - v_i(0, \hat{R}_i) + U_{in} - v_i(0, \hat{R}_i)}. \end{aligned}$$

■

The following lemma will be useful in the proof of Proposition 4.

LEMMA 2. *Let any $m_i, m'_i \in S_i(R^*)$ and x_i . Then*

$$g_i(m_i, x_{-i}) - g_i(m'_i, x_{-i}) \geq -2KX_i.$$

Proof. Let \hat{x}_{-i} such that $\hat{x}_j^{m_j} = x_j^{m_j}$ for all $m_j \in S_j(R^*) \setminus m_j^*$, $\hat{x}_j^{m_j} = 0$ for all $m_j \in \bar{S}_j(R^*)$, and $\hat{x}_j^{m_j^*} = x_j^{m_j^*} + \sum_{m_j \in \bar{S}_j(R^*)} x_j^{m_j}$.

Since $u_i(m_i, x_{-i}) = u_i(m'_i, x_{-i})$ for all $x_{-i} \in S_{-i}(R^*)$, then $g_i(m_i, \hat{x}_{-i}) = g_i(m'_i, \hat{x}_{-i})$.

By Lipschitz continuity we have that

$$g_i(m_i, x_{-i}) - g_i(m_i, \hat{x}_{-i}) \geq -K|x_{-i} - \hat{x}_{-i}| \quad (11)$$

$$g_i(m'_i, \hat{x}_{-i}) - g_i(m'_i, x_{-i}) \geq -K|x_{-i} - \hat{x}_{-i}|. \quad (12)$$

Since $g_i(m_i, \hat{x}_{-1}) = g_i(m'_i, \hat{x}_{-i})$ and $|x_{-i} - \hat{x}_{-i}| = X_i$, the result follows by adding up inequalities (11) and (12). ■

Proof of Proposition 4. By contradiction.

Suppose that (a) is the statement that stops being true earliest, that it does so for agent i and strategy $m_i \in \bar{S}_i(R^*)$ and that the boundary time is t' . Then it must be

$$\frac{x_i^{m_i}(t')}{x_i^{m_i}(0)} = \exp[-h_g t'] \frac{H}{x_i^{m_i^*}(0)}.$$

Notice that, for all t ,

$$\begin{aligned} & u_i(m_i, x_{-i}(t)) - u_i(m_i^*, x_{-i}(t)) \\ & \leq v_i(0, \hat{R}_i) \Pi_{j \neq i} x_j^{m_j^*}(t) + U_i(1 - \Pi_{j \neq i} x_j^{m_j^*}(t)) \\ & \quad - (v_i(f_i(\phi(i, R^*)), \hat{R}_i) \Pi_{j \neq i} x_j^{m_j^*}(t) \\ & \quad + v_i(0, \hat{R}_i)(1 - \Pi_{j \neq i} x_j^{m_j^*}(t))) \\ & = U_i - v_i(0, \hat{R}_i) \\ & \quad - \Pi_{j \neq i} x_j^{m_j^*}(t) (v_i(f_i(\phi(i, R^*)), \hat{R}_i) + U_i - 2v_i(0, \hat{R}_i)). \end{aligned}$$

Since (b) is true for $t < t'$,

$$\begin{aligned} & u_i(m_i, x_{-i}(t)) - u_i(m_i^*, x_{-i}(t)) \\ & < U_i - v_i(0, \hat{R}_i) - H^{n-1} (v_i(f_i(\phi(i, R^*)), \hat{R}_i) + U_i - 2v_i(0, \hat{R}_i)). \end{aligned}$$

Thus,

$$u_i(m_i, x_{-i}(t)) - u_i(m_i^*, x_{-i}(t)) < -h_v,$$

which, by Assumption d.2 and the definition of h_v and h_g , implies that

$$g_i(m_i, x_{-i}(t)) - g_i(m_i^*, x_{-i}(t)) < -h_g.$$

Given $x_i^{m_i^*}(t') \leq H$, if we integrate $g_i(m_i, x_{-i}(t)) - g_i(m_i^*, x_{-i}(t))$ from 0 to t' , we obtain the following:

$$\frac{x_i^{m_i}(t')}{x_i^{m_i}(0)} < \exp[-h_g t'] \frac{H}{x_i^{m_i^*}(0)}.$$

This is a contradiction.

Suppose that (b) is the statement that stops being true earliest, that it does so for agent i , and that the boundary time is t' . Then, it must be true that $x_i^{m_i^*}(t') = H$.

Notice that Lemma 2 implies that, for all $m_i \in S_i(R^*) \setminus \{m_i^*\}$,

$$g_i(m_i^*, x_{-i}(t)) - g_i(m_i, x_{-i}(t)) > -2KX_i(t). \quad (13)$$

Since (a) holds for $t < t'$, (13) implies that

$$\begin{aligned} g_i(m_i^*, x_{-i}(t)) - g_i(m_i, x_{-i}(t)) &> -2K \left(\exp[-2h_g t] \frac{H^2}{(x_i^{m_i^*}(0))^2} X_i(0) \right) \\ &\geq -2K \left(\exp[-2h_g t] \frac{H^2}{(x_i^{m_i^*}(0))^2} X(0) \right). \end{aligned}$$

By integration,

$$\frac{x_i^{m_i^*}(t')}{x_i^{m_i}(t')} \frac{x_i^{m_i}(0)}{x_i^{m_i^*}(0)} > \exp \left[\frac{-2KX(0)}{2h_g} \frac{H^2}{(x_i^{m_i^*}(0))^2} \right] \geq L.$$

Adding over all strategies in $S_i(R^*)$,

$$\frac{x_i^{m_i^*}(t')}{x_i^{m_i^*}(0)} > \frac{x_i^{S_i(R^*)}(t')}{x_i^{S_i(R^*)}(0)} L = \frac{1 - x_i^{\bar{S}_i(R^*)}(t')}{1 - x_i^{\bar{S}_i(R^*)}(0)} L \geq L.$$

This implies $x_i^{m_i^*}(t') > H$ (using the assumption $x_i^{m_i^*}(0)L > H$), which is a contradiction.

Suppose that (c) is the statement that stops being true earliest, that it does so for agent i , and that the boundary time is t' . Then it must be $x_i^{m_i^*}(t')/x_i^{m_i}(t') = (x_i^{m_i^*}(0)/x_i^{m_i}(0))(1/L)$.

By Lemma 2, for all $m_i \in S_i(R^*) \setminus \{m_i^*\}$,

$$g_i(m_i, x_{-i}(t)) - g_i(m_i^*, x_{-i}(t)) > -2KX_i(t). \quad (14)$$

Since (a) holds for $t < t'$, (14) implies that

$$\begin{aligned} g_i(m_i, x_{-i}(t)) - g_i(m_i^*, x_{-i}(t)) &> -2K \left(\exp[-2h_g t] \frac{H^2 X_i(0)}{(x_i^{m_i^*}(0))^2} \right) \\ &\geq -2K \left(\exp[-2h_g t] \frac{H^2}{(x_i^{m_i^*}(0))^2} X(0) \right). \end{aligned}$$

By integration,

$$\frac{x_i^{m_i}(t')}{x_i^{m_i^*}(t')} \frac{x_i^{m_i^*}(0)}{x_i^{m_i}(0)} > \exp \left[\frac{-2KX_0}{2h_g} \frac{H^2}{(x_i^{m_i^*}(0))^2} \right] \geq L,$$

which implies that

$$\frac{x_i^{m_i^*}(t')}{x_i^{m_i}(t')} < \frac{x_i^{m_i^*}(0)}{x_i^{m_i}(0)} \frac{1}{L},$$

which is a contradiction. Since this exhausts all cases the result follows. ■

Proof of Proposition 6. The proof is constructed as follows. We first characterize the limit of the set of restpoints $\text{RE}(\beta)$, and then analyze the stability properties of each of its elements.

We start by observing that, given $\beta \in (0, 1)$, any restpoint must be completely mixed, and it also must be $x_3 > \beta$, as $\Delta \Pi_3(\cdot)$ is always positive in the interior of the state space Δ (because m_3^0 is a weakly dominated strategy). We also know, by continuity of the vectorfield with respect to λ , that every limiting restpoint of the dynamic, as λ goes to zero, must lie in the set of restpoints of the unperturbed dynamic $\text{RE}(G)$.

First, we analyze the limit set of restpoints under CASE 0. In this case, both players 1 and 2 play their strategy m_i^0 with probability 1, that is $x_i = 0$, for $i = 1, 2$. Setting $\dot{x}_1 = 0$ yields the equation

$$\frac{x_1}{\lambda} = \frac{12(\frac{1}{2} - x_1)}{(1 - x_1)(3 + x_1 - x_3(7 - x_2))} \quad (15)$$

and an analogous expression can be obtained for x_2/λ . Denote by x_3^0 a limiting value in a restpoint, if a limit exists, for x_3 . When the limiting values for x_1 and x_2 are zero we have

$$\lim_{\substack{x_i \rightarrow 0 \\ \lambda \rightarrow 0}} \frac{x_i}{\lambda} = \frac{6}{(3 - 7x_3^0)}. \quad (16)$$

Notice that, in this case, if a restpoint exists, it must be $x_3^0 < \frac{3}{7}$, since $x_i/\lambda > 0$. We set $\dot{x}_3/\lambda = 0$, substitute x_i/λ with the expression in (16), solve for x_3 , and substitute x_i , $i = 1, 2$ and λ by their limiting value of zero. The solutions for x_3^0 take the following form:

$$\hat{x}_3^0 = \frac{1 + 7\beta\sqrt{1 - \beta(46 - 49\beta)}}{10} \quad \text{and}$$

$$\check{x}_3^0 = \frac{1 + 7\beta - \sqrt{1 - \beta(46 - 49\beta)}}{10}.$$

Remember that x_3^0 must be a real, positive number, with $\beta < x_3^0 < \frac{3}{7}$. For the expression under the square root at the numerator to be nonnegative, it must be that $\beta \in [0, \frac{23 - 4\sqrt{30}}{49} \approx 0.0222673]$, which determines the feasible range for both roots. Within this interval of values for β , \hat{x}_3^0 (\check{x}_3^0) is a strictly decreasing (increasing) function of β , which has a minimum and a maximum, whose values are $\frac{15 - 2\sqrt{30}}{35}$ (0) and $\frac{2}{10}$ ($\frac{15 - 2\sqrt{30}}{35}$), respectively. As $\beta \rightarrow \frac{23 - 4\sqrt{30}}{49}$, both solutions converge to $\frac{15 - 2\sqrt{30}}{35}$.

We now deal with the subset of limiting restpoints under CASE 1, i.e., with limiting values for $x_i = 1$ for $i = 1, 2$. The equations corresponding to (15) and (16) are now

$$\frac{(1 - x_1)}{\lambda} = \frac{12(x_1 - \frac{1}{2})}{x_1(7x_3 + x_2(1 - x_3) - 3)} \quad (17)$$

$$\lim_{\substack{x_i \rightarrow 1 \\ \lambda \rightarrow 0}} \frac{(1 - x_i)}{\lambda} = \frac{3}{2(2x_3^1 - 1)}, \quad (18)$$

where x_3^1 denotes a limiting value for x_3 (if a limit exists). By analogy with CASE 0, we know from (18) that, if a restpoint exists, it must be $x_3^1 > \frac{1}{2}$. There is a unique feasible solution for x_3^1 , $\forall \beta \in (0, 1)$ which has the following form:

$$\hat{x}_3^1 = \frac{3 + 4\beta + \sqrt{9 - 16\beta(1 - \beta)}}{10}.$$

Following the same procedure for the remaining restpoints of the unperturbed dynamics (i.e., the pure strategy profiles which belong to $\text{RE}(G)$ and do not satisfy either CASE 0 or CASE 1) does not add any element to the limiting set of restpoints of the perturbed dynamics. This should not be surprising, as any other restpoint of the unperturbed replicator dynamics is

unstable with respect to the interior. Since this exhausts all cases, the result follows.

We now move to establish the stability properties of each limiting restpoint separately. The Jacobian matrix $J(x, \lambda)$ for the dynamic system is as follows:

$$\begin{array}{ccc} (1 - 2x_1) \Delta \Pi_1 - \lambda & \frac{-(1 - x_1)x_1(1 - x_3)}{12} & \frac{(1 - x_1)x_1(7 + x_2)}{12} \\ \frac{-(1 - x_2)x_2(1 - x_3)}{12} & (1 - 2x_2) \Delta \Pi_2 - \lambda & \frac{(1 - x_2)x_2(7 + x_1)}{12} \\ \frac{(1 - 2x_2)(1 - x_3)x_3}{6} & \frac{(1 - 2x_1)(1 - x_3)x_3}{6} & (1 - 2x_3) \Delta \Pi_3 - \lambda \end{array}.$$

We analyze CASE 0 first. We know that, in this case, we have two restpoints, which we call $\hat{x}^0 \equiv (0, 0, \hat{x}_3^0)$ and $\check{x}^0 \equiv (0, 0, \check{x}_3^0)$. We evaluate the Jacobian when x_1, x_2 , and λ are equal to their limiting value (i.e., zero). The corresponding eigenvalues are $\{0, (-3 + 7x_3^0)/12, (-3 + 7x_3^0)/12\}$. There are then two (identical) negative eigenvalues (since any limiting $x_3^0 < \frac{3}{7}$ for CASE 0), while the third eigenvalue is equal to zero. To determine the stability properties of the perturbed system, the sign of the eigenvalue whose limit is zero becomes crucial given that continuity of $J(\cdot)$ ensures that the other two will be negative, for any λ sufficiently small. We now linearize the restpoints (as a function of λ) around NE^0 . We set $\check{x}(\lambda, \delta) \equiv (\delta_1 \lambda, \delta_2 \lambda, x_3^0 + \delta_3 \lambda)$, where $\delta \equiv (\delta_1, \delta_2, \delta_3)$ denotes the vector collecting the coefficients of the linearized system. We then evaluate the following expression:

$$\phi^0(x_3^0, \delta) \equiv \lim_{\lambda \rightarrow 0} \frac{\partial \det(J(x, \lambda) |_{\check{x}(\lambda, \delta)})}{\partial \lambda}.$$

We do so because $\det(J(x, \lambda))$, which is equal to zero $\forall x \in NE^0$, will preserve the sign of the third eigenvalue, given that the sign of the other two will stay constant (and negative) when x is sufficiently close to NE^0 and λ is sufficiently small. For CASE 0 we get the following result:

$$\phi^0(x_3^0, \delta) = \frac{-54 + x_3^0(252 + 294x_3^0) + (\delta_1 + \delta_2) \times (9 - 39x_3^0 + 63(x_3^0)^2 - 49(x_3^0)^3)}{864}. \quad (19)$$

We first notice that (19) does *not* depend on δ_3 . To evaluate $\text{sign}(\phi^0(x_3^0, \delta))$ we only need to get estimates of δ_1 and δ_2 , the linear coefficients which measure the responsiveness of the equilibrium values of x_i , $i = 1, 2$ to small changes in λ . We do so by setting $\lim_{\lambda \rightarrow 0} \frac{d}{d\lambda} \tilde{D}(x, \lambda)|_{\tilde{x}(\lambda, \delta)} = 0$ and solving for $\{\delta_1, \delta_2, x_3^0\}$. There are two alternative set of solutions; each of them corresponds to each of the restpoints. In particular,

$$\check{\delta}_1^0 = \check{\delta}_2^0 = \frac{23 - 49\beta - 7\sqrt{1 - \beta(46 - 49\beta)}}{8}$$

$$\hat{\delta}_1^0 = \hat{\delta}_2^0 = \frac{23 - 49\beta + 7\sqrt{1 - \beta(46 - 49\beta)}}{8}.$$

We evaluate the numerator of (19) for both sets of solutions, obtaining the expressions

$$\check{\phi}(\beta) = \frac{3(-7 + 322\beta - 343\beta^2 + (49\beta - 23)\sqrt{1 - 46\beta + 49\beta^2})}{10} \quad (20)$$

$$\hat{\phi}(\beta) = \frac{2863 - 147476\beta + 882882\beta^2 - 1546244\beta^3 + 823543\beta^4 + k\sqrt{146\beta + 49\beta^2}}{1000}, \quad (21)$$

with $k = (3887 - 60123\beta + 165669\beta^2 - 117649\beta^3)$.

Both $\check{\phi}^0(\beta)$ and $\hat{\phi}^0(\beta)$ are plotted in Fig. 4. As the diagram shows, $\check{\phi}^0(\beta)$ is always negative in the domain $[0, \frac{23 - 4\sqrt{30}}{49}]$, whereas $\hat{\phi}^0(\beta)$ is not. In consequence, \check{x}^0 is asymptotically stable whereas \hat{x}^0 is not.

We now move on to CASE 1. Here we have a unique restpoint, which we call $\hat{x}^1 \equiv (1, 1, \hat{x}_3^1)$. The eigenvalues of the unperturbed dynamics are as follows: $\{0, (1 - 2x_3)/3, (1 - 2x_3)/3\}$. As in CASE 0, there are two (identical) negative eigenvalues (given that $x_3 > \frac{1}{2}$), and the remaining eigenvalue is equal to zero. By analogy with CASE 0, we define $\tilde{x}(\lambda) \equiv (1 - \delta_1\lambda, 1 - \delta_2\lambda, x_3^0 + \delta_3\lambda)$ and solve $\lim_{\lambda \rightarrow 0} \frac{d}{d\lambda} \tilde{D}(x, \lambda)|_{\tilde{x}(\lambda, \delta)} = 0$ to get estimates of δ . The unique feasible solution (corresponding to the unique limiting equilibrium) takes the following form:

$$\hat{\delta}_1^1 = \hat{\delta}_2^1 = \frac{3(2 - 4\beta_3 + \sqrt{9 - 16\beta + 16\beta^2})}{2}.$$

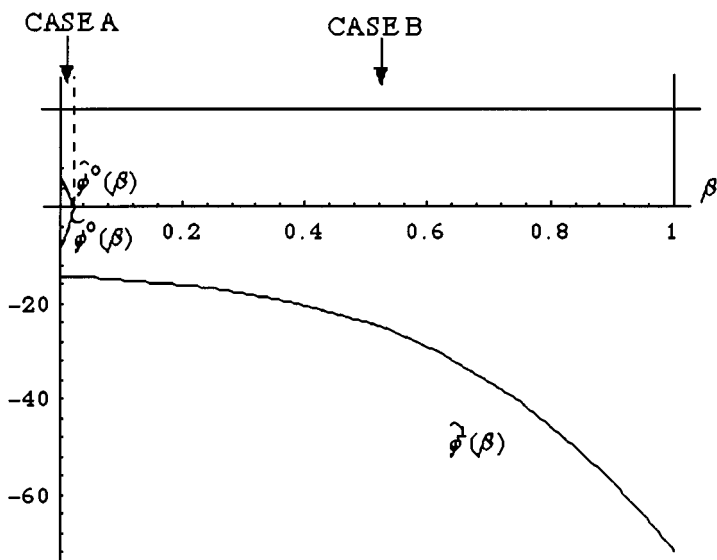


FIG. 4. Asymptotic stability of the dynamic with drift.

The function corresponding to (20) and (21) now takes the form

$$\hat{\phi}^1(\beta) = \frac{24(-\alpha + (2 - 4\beta)\sqrt{\alpha})}{5},$$

with $\alpha = 9 - 16\beta$. The function $\hat{\phi}^1(\beta)$ is also plotted in Fig. 4. As the diagram shows, $\hat{\phi}^1(\beta)$ stays negative $\forall \beta \in (0, 1)$. Thus, \hat{x}^1 is asymptotically stable under any drift configuration. ■

Proof of Proposition 7. For any given arbitrary preference profiles $R \in \Phi$, with $R \neq \hat{R}$, $m_i = \{m_i \notin \hat{S}_i \mid R_{i-1}^i = R_{i-1}, R_{i+1}^i = R_{i+1}\}$ is weakly dominated by $\hat{m}_i = \{m_i \in \hat{S}_i \mid R_{i-1}^i = R_{i-1}, R_{i+1}^i = R_{i+1}\}$, which in turn implies that, for any interior solution $x(t, x(0))$ of (9), $\dot{x}_i^{m_i}(t) = -x_i^{m_i}(t)$ and, therefore,

$$\lim_{t \rightarrow \infty} x_i^{m_i}(t) = 0 \quad (22)$$

for any $m_i \notin \hat{S}_i$. Let $\hat{\Delta}$ denote the face of Δ spanned by the restricted game $(\Gamma, \hat{R})|_{\times \hat{S}_i}$, with $\hat{B}(\varepsilon) = \{x \in \Delta : \|x - \hat{\Delta}\| \leq \varepsilon\}$. An implication of (22) is that $\hat{\Delta}$ is globally interior attracting for the best-reply dynamics (9), as it contains the set of undominated mixed strategies. Furthermore, for all

$i \in I$, \hat{s}_i is a strictly dominant strategy for game $(\Gamma, \hat{R})|_{\times \hat{s}_i}$. This implies that there must exist some positive constants $\bar{\varepsilon}$ and $T(\bar{\varepsilon})$ such that $\text{BR}(x) = \hat{s}$ for any $x \in \hat{B}(\bar{\varepsilon})$ and $x(t) \in \hat{B}(\bar{\varepsilon})$ for any $t \geq T(\bar{\varepsilon})$. We can evaluate $T(\bar{\varepsilon})$ explicitly solving $(1 - \bar{\varepsilon})\exp[-T(\bar{\varepsilon})] = \bar{\varepsilon}$:

$$T(\bar{\varepsilon}) = \ln \left[\frac{1}{\bar{\varepsilon}} - 1 \right]. \quad (23)$$

By virtue of (23), $T(\bar{\varepsilon}) < \infty$. Therefore, the system of differential equations

$$\begin{aligned} \dot{x}_i^{m_i}(t) &= 1 - x_i^{m_i}(t), & m_i &= \hat{s}_i \\ \dot{x}_i^{m_i}(t) &= -x_i^{m_i}(t), & m_i &\neq \hat{s}_i \end{aligned} \quad (24)$$

defines the unique interior solution of (9) for t sufficiently large. This, in turn, implies $\lim_{t \rightarrow \infty} x(t) = \hat{s}$. ■

REFERENCES

- Aubin, J. P., and Cellina, A. (1984). *Differential Inclusions*. Berlin: Springer-Verlag.
- Binmore, K., and Samuelson, L. (1999). "Evolutionary Drift and Equilibrium Selection," *Review of Economic Studies* **66**, 363–394.
- Börgers, T. (1994). "Weak Dominance and Approximate Common Knowledge," *Journal of Economic Theory* **64**, 265–276.
- Börgers, T., and Sarin, R. (1997). "Learning through Reinforcement and Replicator Dynamics," *Journal of Economic Theory* **77**, 1–14.
- Brown, G. W. (1951). "Iterative Solutions of Games by Fictitious Play," in *Activity Analysis of Production and Allocation*. New York: Wiley.
- Cabrales, A. (1999). "Adaptive Dynamics and the Implementation Problem with Complete Information," *Journal of Economic Theory* **86**, 159–184.
- Cabrales, A., Charness, G., and Corchón, L. C. (1998). "An Experiment on Nash Implementation," Universitat Pompeu Fabra Working Paper 300.
- Cabrales, A., and Ponti, G. (1998). "Implementation, Elimination of Weakly Dominated Strategies and Evolutionary Dynamics," UCSB Discussion Paper 6-98.
- Chen, H., Friedman, J., and Thisse, J. F. (1997). "Boundedly Rational Nash Equilibrium: A Probabilistic Choice Approach," *Games and Economic Behavior* **18**, 32–54.
- Chen, Y., and Tang, F. F. (1998). "Learning and Incentive Compatible Mechanisms for Public Goods Provision: An Experimental Study," *Journal of Political Economy* **106**, 633–662.
- Cooper, R., DeJong, D. V., Forsythe, R., and Ross, T. W. (1996). "Cooperation without Reputation: Experimental Evidence from Prisoner's Dilemma Games," *Games and Economic Behavior* **12**, 187–218.
- Cressman, R. (1996). "Evolutionary Stability in the Finitely Repeated Prisoner's Dilemma," *Journal of Economic Theory* **68**, 234–248.

- Cressman, R., and Schlag, K. (1998). "The Dynamic (In)stability of Backward Induction," *Journal of Economic Theory* **83**, 260–285.
- Dekel, E., and Fudenberg, D. (1990). "Rational Behavior with Payoff Uncertainty," *Economic Theory* **52**, 243–267.
- Fudenberg, D., and Kreps, D. (1993). "Learning Mixed Equilibria," *Games and Economic Behavior* **5**, 320–367.
- Gale, J., Binmore, K., and Samuelson, L. (1995). "Learning to Be Imperfect: The Ultimatum Game," *Games and Economic Behavior* **8**, 56–90.
- Gantner, A., Montgomery, R., and Ponti, G. (1998). "Solomon's Dilemma: An Experimental Study on Dynamic Implementation," mimeo, UCSB.
- Groves, T., and Ledyard, J. (1977). "Optimal Allocation of Public Goods: A Solution to the Free Rider Problem," *Econometrica* **45**, 783–809.
- Güth, W., Schmittberger, R., and Schwarze, B. (1982). "An Experimental Analysis of Ultimatum Bargaining," *Journal of Economic Behavior and Organization* **3**, 367–388.
- Hofbauer, J. (1997). "Stability of the Best-Response Dynamics," mimeo, University of Vienna.
- Jackson, M. O., Palfrey, T. R., and Srivastava, S. (1994). "Undominated Nash Implementation in Bounded Mechanisms," *Games and Economic Behavior* **6**, 474–501.
- Jordan, J. S. (1986). "Instability in the Implementation of Walrasian Allocations," *Journal of Economic Theory* **39**, 301–328.
- Marx, L., and Swinkels, J. (1997). "Order Independence for Iterated Weak Dominance," *Games and Economic Behavior* **18**, 219–245.
- Maskin, E. (1977). "Nash Implementation and Welfare Optimality," mimeo, Massachusetts Institute of Technology.
- Matsui, A. (1992). "Best Response Dynamics and Socially Stable Strategies," *Journal of Economic Theory* **57**, 343–362.
- McKelvey, R. D., and Palfrey, T. B. (1992). "An Experimental Analysis of the Centipede Game," *Econometrica* **60**, 803–836.
- Mookherjee, D., and Sopher, B. (1997). "Learning and Decision Costs in Experimental Constant Sum Games," *Games and Economic Behavior* **19**, 97–132.
- Muench, T., and Walker, M. (1984). "Are Groves–Ledyard Equilibria Attainable?," *Review of Economic Studies* **50**, 393–396.
- Nachbar, J. H. (1990). "Evolutionary Selection Dynamics in Games: Convergence and Limit Properties," *International Journal of Game Theory* **19**, 59–89.
- Ponti, G. (2000). "Cycles of Learning in the Centipede Game," *Games and Economic Behavior* **30**, 115–141.
- Repullo, R. (1987). "A Simple Proof of Maskin's Theorem on Nash Implementation," *Soc. Choice Welf.* **4**, 39–41.
- Roth, A., and Erev, I. (1995). "Learning in Extensive Games: Experimental Data and Simple Dynamic Models in the Intermediate Term," *Games and Economic Behavior* **8**, 164–212.
- Samuelson, L. (1993). "Does Evolution Eliminate Weakly Dominates Strategies?," in *Frontiers of Game Theory* (K. Binmore et al., Eds.). London: MIT Press.
- Samuelson, L., and Zhang, J. (1992). "Evolutionary Stability in Asymmetric Games," *Journal of Economic Theory* **57**, 363–391.
- Schlag, K. (1998). "Why Imitate, and If So, How? A Bounded Rational Approach to Multi-Armed Bandits," *Journal of Economic Theory* **78**, 130–156.

- Sjöström, T. (1994). "Implementation in Undominated Nash Equilibria without Integer Games," *Games and Economic Behavior* **6**, 502–511.
- Tang, F. F. (1996). "Anticipatory Learning in Two-Person Games: An Experimental Study: Part II: Learning," Discussion Paper No. B-362, SFB 303, Universität Bonn.
- Taylor, P. D., and Jonker, L. B. (1908). "Evolutionary Stable Strategies and Game Dynamics," *Math. Biosci.* **40**, 145–156.
- Trenquallye, P. de (1988). "Stability of the Groves and Ledyard Mechanism," *Journal of Economic Theory* **46**, 164–171.
- Trenquallye, P. de (1989). "Stable Implementation of Lindahl Allocations," *Economic Letters* **29**, 291–294.
- Vega-Redondo, F. (1989). "Implementation of Lindahl Equilibrium: An Integration of the Static and Dynamic Approaches," *Math. Social Sci.* **18**, 211–228.
- Walker, M. (1981). "A Simple Incentive Compatible Scheme for Attaining Lindahl Allocations," *Econometrica* **49**, 65–71.
- Walker, M. (1984). "A Simple Auctioneerless Mechanism with Walrasian Properties," *Journal of Economic Theory* **32**, 111–127.
- Weibull, J. (1995). *Evolutionary Game Theory*. Cambridge, MA: MIT Press.