

# Implementation, Elimination of Weakly Dominated Strategies and Evolutionary Dynamics

Antonio Cabrales<sup>α</sup>

Universitat Pompeu Fabra and  
ELSE

Giovanni Ponti<sup>γ</sup>

University College London and  
ELSE

December 18, 1996

## Abstract

This paper is concerned with the realism of mechanisms that implement social choice functions in the traditional sense. Will agents actually play the equilibrium assumed by the analysis? As an example, we study convergence and stability properties of Sjstrm's (1994) mechanism, on the assumption that boundedly rational players find their way to equilibrium using monotonic learning dynamics. This mechanism implements most social choice functions in economic environments using as a solution concept the iterated elimination of weakly dominated strategies (only one round of deletion of weakly dominated strategies is needed). There are, however, many sets of Nash equilibria whose payoffs may be very different from those desired by the social choice function. We show that many equilibria in all the sets of equilibria we describe are the limit points of trajectories that have completely mixed initial conditions. The initial conditions that lead to these equilibria need not be very close to the limiting point. Furthermore, even if the dynamics converge to the "right" set of equilibria, it still can converge to quite a poor outcome in welfare terms.

---

<sup>α</sup>This paper was written while Cabrales was visiting the Centre for Economic Learning and Social Evolution at the UCL. He would like to thank for the hospitality and encouragement. The financial support of the EU HCM program is gratefully acknowledged.

<sup>γ</sup>The financial support of the Commission of the European Communities, under the project "Human Capital Mobility - Fellowship ERBCHBICT941575" is gratefully acknowledged.

# 1 Introduction

The theory of implementation studies the problem of designing decentralized institutions through which certain socially desirable objectives can be achieved. These social arrangements should be able to operate without extensive knowledge by the principal about the agents, and in a variety of environments. The principal should ensure that the rules of the game are respected by the agents, and such rules should be designed so that it is in the best interest of the agents to take actions that lead to the socially desirable outcome, given the environment.

More precisely, a social choice rule is implemented by a mechanism (game-form) 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 for every possible environment.

One of the problems that are faced by the implementation theorist is the choice of a solution concept. This is no trivial matter because for different solution concepts the range of social choice rules that can be implemented varies dramatically. As Moore (1990) points out "choice rules are unlikely to be implementable in dominant strategy equilibrium if the domain is very rich and/or the choice rule is efficient". In the case of single-valued choice rules he also notes that "the move from dominant strategy to Nash may not help at all: only the restricted class of strategy-proof choice may be Nash implementable". If the solution concept is more refined, then the domain of the social choice rule can be much larger. In fact, the social choice rule domains are considerably enlarged for subgame-perfect implementation (Moore and Repullo 1988), and even more so when the solution concept is the iterative deletion of weakly dominated strategies (Abreu and Matsushima 1994, Jackson et al. 1994, Sjåstråm 1994). In fact, as Sjåstråm (1994) says: "With enough ingenuity the planner can implement anything".

The question that arises then is whether the equilibrium concept chosen is a good one for the game in object. One way to answer the question is to assume that agents are boundedly rational and that they adjust their actions over time through some trial and error learning procedure. One can then analyze under which conditions the actions that lead to the socially desirable outcomes are played asymptotically, if at all. Research in implementation theory has paid little attention to the problem of how an equilibrium is reached. Some exceptions are the papers of Muench and Walker (1984) and de Trenqualye (1988), who study the local stability of the Groves and Ledyard (1977) mechanism, and Cabrales (1996), who studies the global convergence of the canonical mechanism (Maskin 1977, Repullo 1987) of Nash implementation and the mechanisms of Abreu and Matsushima (1992, 1994).

In this paper we study the convergence and stability properties of the Sjåstråm's (1994) mechanism when one assumes that the players are boundedly rational and the dynamics are monotonic (Samuelson and Zhang 1992, Weibull 1995). One particularly well known

member of the family of monotonic dynamics is the replicator dynamics of evolutionary game theory (Taylor and Jonker, 1978). These dynamics have been given a learning theoretic foundation by Bargers and Sarin (1993), and they can also be interpreted as a model of imitation (Schlag, 1994). Sjastram's (1994) mechanism and the one that Jackson et al. (1994) study for separable environments are very similar and most of our results would generalize easily for that mechanism as well.

We concentrate on Sjastram's mechanism for several reasons. One is that the conditions for implementation are quite weak. Although the environments that are permitted are not universal, they are rich enough for most economic problems. Furthermore, this reduction in the domain permits the author to implement the social choice rule with a "bounded" game and thus makes it immune to the criticisms of Jackson (1992). Finally, although the solution concept is the iterated elimination of (weakly) dominated strategies (it also implements in undominated Nash equilibria), it only needs one round of deletion of weakly dominated strategies (the "first"). This last feature of the mechanism makes it particularly attractive since under some assumptions of imperfect knowledge of agents (either because of payoff uncertainty as in Dekel and Fudenberg, 1990, or through lack of perfect common knowledge of rationality as in Bargers, 1994) the appropriate solution concept implies one round of deletion of weakly dominated strategies and then the iterated deletion of strictly dominated strategies.

In Sjastram's (1994) mechanism the agents are arranged to announce their preferences and those of their two closest neighbors. The mechanism is designed in such a way that a truthful report of one's own preferences is weakly dominant (it does not affect one's payoff, except in a set of states which is called totally inconsistent, and in those states it is preferable to report them truthfully). Since in this mechanism it is advantageous to report the same preferences about your neighbors that 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. There are, however, many other Nash equilibria. 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 they report the preferences about themselves indicated in  $R$  with high enough (this need not be very high) probability. The reason for this is that the mechanism makes it important that all agents match their neighbors' announcements about themselves, but the report about oneself is only important in some unlikely (totally inconsistent) state.

We show that many equilibria in all the components of equilibria we have described are the limit points of trajectories of the learning dynamics that have completely mixed initial conditions (that is, initial conditions that give strictly positive weights to all possible messages). Although the general results are local, we can show by example (the game in Figure 1, Sjastram, 1994) that the initial conditions that lead to these equilibria need not be close to the limiting point. Furthermore, and perhaps more worrying, the equilibria which belong to the same component as the completely truthful report are not outcome equivalent to such equilibrium, as they yield payoffs that are significantly different (lower)

to the payoffs of the social choice functions outcome. Therefore, even if the dynamics converge to the "right" component of equilibria, it still can converge to quite a poor outcome in welfare terms.

One could naively expect that evolution would eliminate weakly dominated strategies. The reason why this doesn't happen is that the weakly dominant strategy grows faster than the dominated one only if the totally inconsistent states are met often enough by the players. But the weight of the totally inconsistent states is also decreasing over time since people are learning to avoid such states. It may be that they decrease fast enough so that the push towards the weakly dominant strategy is not enough to make the dominated strategy disappear.

The fact that evolution does not eliminate weakly dominated strategies has been known since at least Nachbar (1990). Samuelson (1993) discusses the issue of elimination of weakly dominated strategies in evolutionary games. Binmore et al. (1995) have shown the implications of these findings for the ultimatum bargaining game. In particular, they provide a numerical example, based on the classic "chain store game", in which a) there are trajectories of the replicator dynamics which converge to the Nash equilibrium component in which the players choose a weakly dominated strategy with positive probability and b) in the presence of mutations, such component may even exhibit asymptotic stability properties. These results are more than a theoretical curiosity. Binmore and Samuelson (1996) note that: "the experimental evidence is now strong that one cannot rely on predictions that depend on deleting weakly dominated strategies".

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 the dynamics for the game in Figure 1, Sjöström 1994, to be considered a simplified version of the mechanism. In section 4 we give local results (for some interior initial conditions) for the set of limit points of the dynamics for the general game. In section 5 we describe the asymptotic stability properties of the sets of limit points in the presence of mutations. Finally, section 6 concludes, together with an appendix containing the proofs of the relevant propositions.

## 2 The model and the dynamics

We will first introduce some notation, the assumptions on preferences and the game form proposed by Sjöström (1994). Then we introduce the assumptions we make on the dynamics.

The only important change we make in the presentation with respect to Sjöström (1994) is that we make the assumptions on a Von Neumann-Morgenstern utility function instead of on the preference relation. We do this because we need to specify the payoff functions

for mixed strategies, since the dynamics are defined on the mixed strategy simplex.

There is a set  $I = \{1, \dots, n\}$  (with  $n \geq 3$ ) of agents and a set  $A \subseteq \mathbb{R}_+^m$  of feasible consumption plans. The preferences of agent  $i \in I$  are represented with a (Von Neumann-Morgenstern) utility function  $v_i : A \rightarrow \mathbb{R}$ , where  $\mathbb{R}$  specifies a finite set of possible utility functions. An element  $R_i$  of  $\mathbb{R}$  is meant to represent the preferences of agent  $i$  over  $A$ . A preference profile is a vector  $R = (R_1, \dots, R_n)$ , where  $R_i \in \mathbb{R}$ . The preference profiles will be common knowledge among the agents. The following assumptions are made concerning preferences and feasible consumption profiles.

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

**Assumption p.2.** The preferences represented by  $R_i \in \mathbb{R}$  are complete and transitive.

**Assumption p.3.** The preferences represented by  $R_i \in \mathbb{R}$  are strictly convex. That is, for any  $a, a^0 \in \mathbb{R}_+^m$  and for all  $\lambda \in (0, 1)$  if  $a \succ a^0$  and  $v_i(a; R_i) > v_i(a^0; R_i)$  then  $v_i(\lambda a + (1 - \lambda)a^0; R_i) > v_i(a^0; R_i)$ :

**Assumption p.4.** For any  $R_i \in \mathbb{R}$  if  $a \succeq 0$  and  $a \neq 0$  then  $v_i(a; R_i) > v_i(0; R_i)$ :

**Assumption p.5.** Preference reversal. For any  $R_i, R_i^0 \in \mathbb{R}$  if  $R_i \neq R_i^0$  then there are  $a, a^0 \in A$  such that  $v_i(a; R_i) > v_i(a; R_i^0)$  and  $v_i(a^0; R_i^0) > v_i(a^0; R_i)$ :

For any set  $B \subseteq \mathbb{R}_+^m$  and any  $R_i \in \mathbb{R}$  a choice correspondence is defined as follows:  $c(B; R_i) = \{a \in B \mid v_i(a; R_i) \geq v_i(b; R_i) \forall b \in B\}$ :

A social choice function is a mapping  $f : \mathbb{R} \rightarrow A$ ; where  $f(R) = (f_1(R), \dots, f_n(R))$ . A mechanism is a pair  $\Gamma = (M; \Phi)$ , where  $M = \prod_{i \in I} M_i$  and  $\Phi : M \rightarrow A$ ; where  $\Phi(m) = (\Phi_1(m), \dots, \Phi_n(m))$ .  $M_i$  is the message space of agent  $i$  and  $\Phi$  is the outcome function. A mechanism and a preference profile define a game.

Let  $M_{-i} = \prod_{j \in I, j \neq i} M_j$ . Given a mechanism  $\Gamma$  and a preference profile  $R$ , we say that  $m_i$  is weakly dominated for some set of messages  $F \subseteq \prod_{i \in I} F_i \subseteq M$  if there exists a message  $m_i^0 \in F_i$  such that  $v_i(\Phi_i(m_i^0; m_{-i}); R_i) \geq v_i(\Phi_i(m_i; m_{-i}); R_i)$  for all  $m_{-i} \in F_{-i}$  and there is some  $m_{-i}^a \in F_{-i}$  such that  $v_i(\Phi_i(m_i^0; m_{-i}^a); R_i) > v_i(\Phi_i(m_i; m_{-i}^a); R_i)$ : Define the set  $U_i(F; (i; R)) = \{m_i \in F_i \mid m_i \text{ is not weakly dominated in } F \text{ for the game } (i; R)\}$ :

The message  $m_i$  is a best response for player  $i$ , to  $m_{-i} \in M_{-i}$  if  $v_i(\Phi_i(m_i; m_{-i}); R_i) \geq v_i(\Phi_i(m_i^0; m_{-i}); R_i)$  for all  $m_i^0 \in M_i$ . A message profile  $m$  is a Nash equilibrium (NE) if  $m_i$  is a best response to  $m_{-i}$  for all  $i \in I$ . A message profile  $m \in M$  is an undominated Nash equilibrium (UNE) for the game  $(i; R)$  if it is a Nash equilibrium and  $m_i \in U_i(M; (i; R))$ : Let  $UNE(i; R) = \{f^{\Phi}(m) \in A \mid m \text{ is an UNE for the game } (i; R)\}$ :

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

For the iterated deletion of weakly dominated strategies let  $U_i^1(i; R) = U_i(M; (i; R))$ , and if  $U_i^k(i; R)$  has been defined for  $k \leq 1$ , let  $U_i^{k+1}(i; R) \hat{=} U_i(\bigcup_{j \neq i} U_j^k(i; R); (i; R))$ : Let  $U_i^1(i; R) \hat{=} \bigcap_{k=1}^{\infty} U_i^k(i; R)$ : Let  $IWD(i; R) \hat{=} f^{\otimes}(m) \geq A_j m_i \geq U_i^1(i; R)$  for all  $i$ g:

We say that a mechanism  $(M; g)$  implements a social choice function  $f$  with iterated deletion of weakly dominated strategies if for all  $R \in \mathcal{C}$ ,  $f(R) = IWD(R)$ :

We now construct a mechanism.

Let  $M_i = \mathcal{C}_{i-1} \times \mathcal{C}_i \times \mathcal{C}_{i+1}$ , 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  $f m_{j_1}; \dots; m_{j_K} g$  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  $f m_{j_1}; \dots; m_{j_K} g$  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  $R_i; R_i^0 \in \mathcal{C}_i$ , where  $R_i \neq R_i^0$ . By assumption p.6 there are  $a; a^0 \in A$  such that  $v_i(a; R_i) > v_i(a; R_i^0)$  and  $v_i(a; R_i^0) > v_i(a^0; R_i^0)$ : We can choose  $a$  and  $a^0$  so that  $v_i(a; R_i) > v_i(a^0; R_i)$  for all  $a^0$  in the line segment between  $a$  and  $a^0$ : Given this pair  $(a; a^0)$  let  $\bar{c}_i(R_i; R_i^0) \hat{=} f b \in \mathcal{C}_i^+ | b = \lambda a + (1 - \lambda)a^0; \text{ for } \lambda \in [0; 1] g$ : By construction, for all  $R_i; R_i^0 \in \mathcal{C}_i$ ;  $\bar{c}_i(R_i; R_i^0); R_i \succ \bar{c}_i(R_i; R_i^0); R_i^0$ : Let  $\hat{A}(i; m) \hat{=} (R_1^i; R_2^i; \dots; R_{i-1}^i; R_{i+1}^i; R_{i+2}^i; \dots; R_n^i)$  and for every  $i$  and  $m_i$ ; define

$$B_i(m_i) = \begin{cases} \bar{c}_i(\hat{A}(i; m)) & \text{if } m_i \text{ is totally consistent} \\ \bar{c}_i(R_{i-1}^i; R_{i+1}^i) & \text{if } m_i \text{ is totally inconsistent} \\ \frac{1}{n} f_i(\hat{A}(i; m)) & \text{Otherwise} \end{cases}$$

Now we can define  $\mathcal{C}^{\otimes}$ :

$$\mathcal{C}^{\otimes}(m) = \begin{cases} \bar{c}(B_i(m_i); R_i^i) & \text{if } R_{i-1}^i = R_{i-2}^i \text{ and } R_{i+1}^i; R_{i+2}^i \\ 0 & \text{otherwise} \end{cases}$$

To understand the mechanism notice that the only time when the choice of an announcement  $R_i^i$  has any effect on payoffs is when  $m_i$  is totally inconsistent. In that case, the outcome is the optimal choice within the set  $\bar{c}_i(R_{i-1}^i; R_{i+1}^i)$  according to the announced  $R_i^i$ . For this reason announcing the true preference  $R_i^i$  can never hurt. Furthermore, for every alternative announcement  $R_i^i = \hat{R}_i$ , there is some totally inconsistent  $m_i$  with  $R_{i-1}^i = R_{i-1}^i$  and  $R_{i+1}^i = \hat{R}_i$  and the set  $\bar{c}_i(\cdot; \cdot)$  is constructed in such a way that  $\bar{c}_i(R_i^i; \hat{R}_i); R_i^i$  is strictly preferred to  $\bar{c}_i(R_i^i; \hat{R}_i); \hat{R}_i$ . Therefore, a message  $m_i = (R_{i-1}^i; \hat{R}_i; R_{i+1}^i)$  is weakly dominated by a message  $m_i = (R_{i-1}^i; R_i^i; R_{i+1}^i)$ , that is, untruthful announcements about oneself are weakly dominated.

Once these weakly dominated strategies are eliminated and all agents announce the true preferences about themselves,  $R_j^i = R_j^i$ , it is strictly dominated to announce untruthful preferences about the neighbors,  $R_{i+1}^i \neq R_{i+1}^i = R_{i+1}^i$  or  $R_{i-1}^i \neq R_{i-1}^i = R_{i-1}^i$ , since disagreeing with the neighbors is punished with the 0 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 will undertake below, that the set of states for which not announcing the true preferences about oneself is weakly dominated are themselves states that typically produce very bad outcomes for the opponents (at least one of them will have 0 consumption, and probably many). If agents learn to avoid totally inconsistent states very fast, there is no incentive to tell the truth about oneself. The mechanism we have described puts a lot of emphasis in consensus announcements, since disagreement is punished with 0 consumption, and truth-telling is only rewarded in a set of states which need not be very prominent in the minds of the players. That is precisely the reason why convergence to outcomes of the social choice function may fail to occur. This conflict is typical of other mechanisms that implement in the iterated deletion of weakly dominated strategies, like Abreu and Matsushima (1994).

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

Fix a given preference profile  $R \in \mathcal{R}^I$ : Let  $x_i^{m_i}$  be the probability assigned by agent  $i$  to message  $m_i$ ; and  $x_i \in \Delta_i$  be a mixed strategy for agent  $i$  (where  $\Delta_i$  denotes the  $|M_i|$ -dimensional simplex which describes player  $i$ 's mixed strategy space). Let also  $x_{-i} \in \prod_{j \neq i} \Delta_j$  be a mixed strategy profile for agents other than  $i$ ; with  $x = (x_i; x_{-i}) \in \Delta = \prod_{i \in I} \Delta_i$ . Finally, let  $u_i(x_i; x_{-i}) = \sum_{m_i \in M_i} v_i(m_i; m_{-i}; R_i) x_i^{m_i}$ :

We formalize player  $i$ 's behavior in terms of the mixed strategy  $x_i(t)$  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  $\Delta$ , with  $\Delta^0$  denoting the relative interior of  $\Delta$ , that is, the set of completely mixed strategy profiles. We make the following assumption:

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

$$\dot{x} = D(x(t)) \quad (1)$$

We require that the autonomous system (1) satisfies the standard regularity condition, i.e.,  $D$  must be i) Lipschitz continuous with ii)  $\lim_{k \rightarrow \infty} D_i^k(\cdot) = 0$  and iii)  $\lim_{x_i^{m_i} \rightarrow 0} \frac{x_i^{m_i}}{x_i^{m_i}}$  well-defined and finite.<sup>1</sup> Furthermore,  $D$  must also satisfy the following requirements:

**Assumption d.2.**  $D$  is a regular (payoff) monotonic selection dynamic. More explicitly, let  $g_i(m_i; x_{-i}(t)) = \frac{x_i^{m_i}(t)}{x_i^{m_i}(t)}$ . Then for all  $m_i; m_i^0$ ; we have that

$$\text{sign}[g_i(m_i; x_{-i}(t)) - g_i(m_i^0; x_{-i}(t))] = \text{sign}[u_i(m_i; x_{-i}(t)) - u_i(m_i^0; x_{-i}(t))]$$

<sup>1</sup>A useful implication of this regularity assumption is that the solution of the dynamical system leaves  $\Delta$ , as well as  $\Delta^0$ , invariant (and, a fortiori, forward invariant): any solution path starting from  $\Delta$  ( $\Delta^0$ ) does not leave  $\Delta$  ( $\Delta^0$ ). This property will prove to be useful to obtain some of the results of the paper.

Assumption d.2 is a common assumption in the evolutionary literature and we will not dwell on it (see, for example, Samuelson and Zhang 1992 and Weibull, 1995).

**Assumption d.3.** Let  $Y_i(m_i; m_i^0) = \int u_i(m_i; x_i(t)) \cdot u_i(m_i^0; x_i(t)) = 0$ . Then, for all  $\epsilon > 1$ :

$$\frac{\lim_{d(x_i(t); Y_i(m_i; m_i^0)) \rightarrow 0} \sup [g_i(m_i; x_i(t)) - g_i(m_i^0; x_i(t))]}{\text{sign}[u_i(m_i^0; x_i(t)) - u_i(m_i; x_i(t))] (\ln |u_i(m_i^0; x_i(t)) - u_i(m_i; x_i(t))|)^{\epsilon}} > \epsilon$$

Assumption d.3 is less standard in the evolutionary literature and we will expand on it when we discuss Proposition 4 because it will be helpful to understand why weakly dominated strategies need not disappear in the limit. What assumption d.3 says is that if the difference in payoffs between two strategies is going to zero a rate  $\exp[-n]$ , the difference in growth rates has to go to zero at least at a rate  $1/n^\epsilon$ . Continuity and assumption d.2 demand that strategies that have the same payoff grow at the same rate, but they impose no requirements on the speed at which the difference in growth rates goes to zero as the difference in payoff go to zero. Assumption d.3 can be satisfied even if the sensitivity of growth rates to payoffs is much higher than linear around zero (as would be implied, for example, by the replicator dynamics and other aggregate monotonic dynamics).

**Assumption d.4.**  $x(0) \in \mathbb{C}^0$

Finally, Assumption d.4, which is also standard in the evolutionary literature, is necessary because it excludes the possibility that the selection dynamic acts only on a subset of the strategy space. This possibility arises because the system is forward invariant, and therefore a strategy that has zero weight at time zero would also have zero weight at all subsequent times. We want to avoid this possibility because the selection dynamics would be operating on a game which might be qualitatively different from the game we are trying to analyze.

### 3 An example.

We preface the analysis of the dynamics 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 a unit of good which has to be divided among three players: 1, 2 and 3. The (Von-Neumann Morgenstern) utility functions of players 1 and 2 are linear in the amount of good consumed, and this is common knowledge among the players and the planner. The utility function of player 3 may have one of two possible types, either linear on the amount of good consumed (we index these preferences by the number 1) or linear until the amount of good consumed is 1/3, for consumptions larger than 1/3 the utility remains constant at the value 1/3, since the agent becomes satiated (the index for



these preferences is 0). The true preferences of player 3 are common knowledge among the players, but the planner does not know them.

The social choice function recommends the consumption vector  $(1=3; 1=3; 1=3)$  for preferences of type 1 and  $(1=4; 1=4; 1=2)$  for preferences of type 0. Notice that this social choice function is such that agent 3 would like to conceal her preferences, and therefore the planner needs a nontrivial mechanism to elicit the true preferences.

The mechanism proposed by Sjöström requires the three players to make a simultaneous statement about the preferences of player 3. Let  $m_i^1; i \in \{1, 2, 3\}$  represent the message in which preferences of type 1 are announced, with  $m_i^0$  denoting the announcement of type 0 preferences. Figure 1 illustrates the outcome function. We will assume for the analysis that the true preferences are of type 1 and therefore Figure 1 is also the payoff function of a game, which we call  $\Gamma$ :

Figure 1

Sjöström's Example: game  $\Gamma$ :

Player 1 picks a row, player 2 a column, and player 3 picks a matrix. We first notice that the mechanism leads to a game which is weakly dominance solvable, in the sense that it can be reduced to a single cell, corresponding to the truth-telling equilibrium outcome, by the iterated deletion of weakly dominated strategies. Unlike other weakly solvable games, this procedure yields, in this example, a unique outcome, independently of the order of removal of strategies. Player 3 first deletes her (weakly) dominated strategy  $m_3^0$  (the other agents have no dominated strategies at this stage). Once  $m_3^0$  is removed, the strategies  $m_1^0$  and  $m_2^0$  for players 1 and 2 become strictly dominated. The unique strategy profile selected in this way is  $(m_1^1; m_2^1; m_3^1)$ . Notice, however, that the strategy profile  $(m_1^0; m_2^0; m_3^1)$  is also an equilibrium, and that this equilibrium yields a higher payoff for agent 3 than  $(m_1^1; m_2^1; m_3^1)$ .

Given that each player has only two strategies in her support, with an abuse of notation we set  $x_i \in \{0, 1\}$ .<sup>2</sup> We first characterize the set of Nash equilibria of the game:

**Proposition 1.** The set NE of Nash equilibria of  $\Gamma$  is the union of precisely two disjoint components  $NE^0$  and  $NE^1$ , where:

$$NE^0 = \{x \in \{0, 1\}^3 \mid x_1 = x_2 = 0; x_3 = \frac{3}{7}g\}$$

$$NE^1 = \{x \in \{0, 1\}^3 \mid x_1 = x_2 = 1; x_3 = \frac{1}{2}g\}$$

**Proof.** See the Appendix.2

<sup>2</sup>The fact that each player has only two available options will also allow us to express the dynamics in terms of the payoff difference between player  $i$ 's truthful and untruthful strategy, which we call  $\Phi_i(x(t))$  (i.e.  $\Phi_i(x(t)) = u_i(m_i^1; x_{-i}(t)) - u_i(m_i^0; x_{-i}(t))$ ).

Denote with  $RE(i)$  the set of restpoints of  $\dot{x}_i$  under any monotonic dynamic. It is straightforward to show that  $RE(i)$  contains (together with all the pure strategy profiles) only the following components:  $RE^0 = \{x \in \Delta \mid x_1 = x_2 = 0; x_3 \in [0; 1]\}$  and  $RE^1 = \{x \in \Delta \mid x_1 = x_2 = 1; x_3 \in [0; 1]\}$ . 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 solution  $x(t; x(0))$  of a monotonic selection dynamics  $\dot{x} = D(x)$  with completely mixed initial conditions converges asymptotically to NE.

**Proof.** See the Appendix 2

If initial conditions are completely mixed, we then know that the evolutionary dynamics will eventually converge to a Nash equilibrium of the game. In the following section we extend the result to the more general setting of Sjöström's mechanism.

## 4 Local results for the general game

In this section we show that the results of the previous section generalize locally. Proposition 3 characterizes some components of Nash equilibria for the game induced by the mechanism in Sjöström (1994), which we described in section 2. Any message profile in which the agents are unanimous in the (arbitrary) preference profile they announce,  $R^a$ , (or, more appropriately, the preferences they announce about their neighbors and themselves are taken from the profile  $R^a$ ) is an equilibrium. Furthermore, a mixed strategy profile in which every agent mixes between messages consistent with  $R^a$  and other preference profiles that only differ in the announcement they make about their own preferences is also an equilibrium, as long as  $R^a$  is given a high enough weight. As we showed in the example, the weight given to  $R^a$  need not be very high. The equilibria in a component are not payoff equivalent, since disagreeing with a neighbor (an event with nonzero probability in the mixed strategy equilibria) results in a punishment. Proposition 4 shows that any of the previous equilibria that gives enough weight to  $R^a$  is the limit point of some interior path for a monotonic selection dynamic. Figure 2 shows that the initial condition need not be very close to the limit point.

Let  $\hat{R}^i$  be the true preference profile and  $R^a$  an arbitrary preference profile. Let  $m_i^a = (R_{i-1}^a; R_i^a; R_{i+1}^a)$ ;  $U_i = \max_R v_i(f_i(R); \hat{R}_i)$  and  $U_i^a = \max_R v_i(\frac{1}{n}f_i(R); \hat{R}_i)$ .  $m_i^a$  is a consensus announcement by agent  $i$ ,  $U_i$  is utility associated to the most preferred outcome from the social choice function for agent  $i$  with true preferences  $\hat{R}_i$  and  $U_i^a$  is utility associated to the most preferred consumption bundle among those that result from dividing the bundles assigned by the social choice function by  $n$ .

$$S_i = \{m_i \in M_i \mid R_{i-1}^i = R_{i-1}^a; R_i^i = R_i^a; R_{i+1}^i = R_{i+1}^a\}$$

and  $\hat{S}_i = \{m_i \in M_i \mid m_i \in S_i\}$ . The set  $S_i$  is the set of all mixed strategies in which

announcements about the neighbors agrees with  $R^a$ , and  $S_i^d$  is the complement of  $S_i$  with respect to  $M_i$ .

$$S_i^{k_i} = \{x_i \mid x_i^{m_i} = 0; \text{ for all } m_i \in S_i \text{ and } x_i^{m_i} > k_i g_i\}$$

where we assume

$$(k_i)^n \leq \frac{U_{j \in i} u_j(0; \hat{R}_j)}{v_j(f_i(\hat{A}(j; R^a)); \hat{R}_j) + v_j(0; \hat{R}_j) + U_{j \in i} v_j(0; \hat{R}_j)}$$

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

**Proposition 3.** For all  $\hat{R}; R^a \in \mathcal{R}$  and  $x_i \in S_i^{k_i}$ ,  $x$  is a Nash equilibrium of  $(\mathcal{G}; \hat{R})$ .

**Proof.** See the Appendix.2

Now we prove that not only are there other Nash equilibria, but that elements in those components can be reached by paths starting in the interior to the simplex. By assumption d.2 we know that for all  $h_v > 0$  with  $u_i(m_i; x_{-i}(t)) - u_i(m_i^0; x_{-i}(t)) < h_v$ , there is  $h_g > 0$ , such that  $g_i(m_i; x_{-i}(t)) - g_i(m_i^0; x_{-i}(t)) < h_g$ :

Let  $h_v$  be a constant such that  $0 < h_v \leq \min_{i \in R} v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) - v_i(0; \hat{R}_i)$ . Let the corresponding  $h_g$  and

$$H = \max_{i \in R} \frac{U_{j \in i} v_j(0; \hat{R}_j) + h_v}{v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) + U_{j \in i} 2v_j(0; \hat{R}_j)} > 0$$

Note that by definition  $H \in [0; 1]$ :

By Assumption d.3 we know that there exists  $\mu_i(m_i; m_i^0) > 0$  such that if  $|u_i(m_i; x_{-i}(t)) - u_i(m_i^0; x_{-i}(t))| < \mu_i(m_i; m_i^0)$

$$\frac{\limsup_{d(x_{-i}(t); Y_i(m_i; m_i^0)) \rightarrow 0} |g_i(m_i; x_{-i}(t)) - g_i(m_i^0; x_{-i}(t))|}{|\ln |u_i(m_i^0; x_{-i}(t)) - u_i(m_i; x_{-i}(t))||} > 1$$

Let  $0 < \mu < \min_{i; m_i; m_i^0} \mu_i(m_i; m_i^0)$ :

For any set  $\mathcal{C}_i \in M_i$ , let  $x_{\mathcal{C}_i} = \prod_{m_i \in \mathcal{C}_i} x_i^{m_i}$ , and  $L = \exp \left( \frac{h_g (\pm 1) - \ln \left( U_{j \in i} x_j^{s_j}(0) \frac{H}{x_i^{m_i}(0)} \right)}{h_g (\pm 1) - \ln \left( U_{j \in i} x_j^{s_j}(0) \frac{H}{x_i^{m_i}(0)} \right)} \right)$

**Proposition 4.** Assume that for all  $i$ ,  $x_i^{m_i}(0)$  is big enough so that,  $x_i^{m_i}(0)L > H$  and  $U_{j \in i} x_j^{s_j}(0) \frac{H}{x_i^{m_i}(0)} < \mu$ . Under these conditions we have that

- a) For all  $m_i \in \hat{S}_i$ ,  $\frac{x_i^{m_i}(t)}{x_i^{m_i^*}(0)} < \exp[-h_g t] \frac{H}{x_i^{m_i^*}(0)}$  for all  $t$  and all  $i$
- b)  $x_i^{m_i^*}(t) > H$  for all  $t$
- c)  $\frac{x_i^{m_i^*}(t)}{x_i^{m_i^*}(0)} < L$  for all  $t$  and all  $m_i \in S_i$

**Proof.** See the Appendix.2

Part a) of the Proposition says that the weight of any strategy in  $\hat{S}_i$  decreases over time at a rate higher than  $h_g$ . This is important because the strategies for which not telling the truth about oneself is dominated are all in  $\hat{S}_j$ , so if the weight of these strategies decrease over time, the payoff advantage of the dominating strategy disappears over time, and makes it possible for a dominated strategy to have positive limiting weight.

Part b) ensures that the weight of  $m_i^*$  is always high enough. If the weight of  $m_i^*$  is high enough, then the strategies in the sets  $\hat{S}_j$  have a lower payoff than strategies in  $S_j$  since an announcement about your neighbor that does not coincide with her announcement about herself is punished.

In fact parts a) and b) reinforce each other. While  $m_i^*$  keeps having a high enough weight, the weight of strategies in  $\hat{S}_j$  decreases, and if strategies in  $\hat{S}_j$  decrease fast enough the weight of  $m_i^*$  does not go below a certain bound. All of this provided that  $m_i^*$  started with high enough weight, which as Figure 2 demonstrates, need not be very high.

Notice that part b) guarantees that pure strategy equilibria in the "wrong" component are attractors of interior paths. Part c) says that the weight of  $m_i^*$  in the limit is less than 1, and therefore some mixed strategy equilibria are attractors as well if the initial conditions give sufficiently little weight to strategies in  $\hat{S}_j$ . This guarantees that even if there is convergence to the "right" component it need not be to the pure strategy equilibrium, and remember that the equilibria are not payoff equivalent (the mixed strategy equilibria have lower expected payoff because agents are punished for announcing discordant preferences).

The convergence to the mixed equilibria can happen because the payoffs to all strategies in  $S_i$  are similar if the weight of strategies in  $\hat{S}_j$  is small, and by a) the weight of strategies in  $\hat{S}_j$  is decreasing. So even though  $m_i^*$  has a payoff advantage, the advantage goes to zero over time, and assumption d.3 guarantee that it does not accumulate fast enough.

If d.3 didn't hold, equilibria which do not implement the social choice function may fail to be a limit point for the dynamics. Convergence to the "wrong" equilibrium obtains only if a weakly dominated strategy for player  $i$  (call it  $m$ , and call  $m^0$  the strategy that weakly dominates  $m$ ) gets positive weight in the limit. But along the way to the limit the strategy against which  $m$  and  $m^0$  differ (call it  $y$ ) has also positive weight (since the system is regular, and therefore forward invariant). So by assumption d.2 the growth rate of  $m^0$  is larger than the growth rate of  $m$ . The weight of  $m^0 = m$  is the integral of the difference in growth rates of  $m^0$  and  $m$ . If the limiting value of this integral is infinite the limiting value of  $x_i^m(t)$  would be zero. But the weight of  $y$  (and thus the difference in

payoffs) may be going to zero, thus the weight of  $m^0 = m$  is an integral of a function that goes to zero, which may be finite.

Assumption d.3 describes how the growth rates have to relate to payoffs (when differences in payoffs are small) so that the limiting value of a dominated strategy is not zero. Assumption d.3 would hold, for example, if the growth rates were linear in the payoffs, as it happens with the replicator dynamics. But the requirement is much weaker than that, because it is only a local requirement around zero, and because the rate at which growth rates go to zero with payoffs can be much higher than linear. In other words, even if the growth rates were much more sensitive to payoff differences (around zero) than the replicator dynamics allow, assumption d.3 could still be satisfied.

The elimination of a weakly dominated strategy in an evolutionary context requires that the strategy against which the dominated strategy gives a lower payoff than the dominating strategy has to appear sufficiently often or that its appearance has to provoke a dramatic enough reduction in the dominated strategy.

## 5 More on the example (stability with/out drift).

In the previous section, we have extended the convergence result contained in Proposition 2, and we have shown 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 the example in order to test the stability properties of NE. To do so, some further terminology is needed:

**Definition 1.** Let  $C$  be a closed set of restpoints in  $\Phi$  of the dynamics (1): Then: (i)  $C$  is called (interior) stable if, for every neighborhood  $O$  of  $C$ , there is another neighborhood  $U$  of  $C$ , with  $U \subset O$ ; such that the trajectory of any initial point in  $U \setminus \Phi$  ( $U \setminus \Phi^0$ ) remain inside  $O$ ;

(ii)  $C$  is called (interior) attracting if is contained in an open set  $O$  such that every initial point in  $O$  and also in (the interior of)  $\Phi$  evolves under  $D$  to  $C$ .  $C$  is globally (interior) attracting if every point in (the interior of)  $\Phi$  evolves under  $D$  to  $C$ ;

(iii)  $C$  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.1-5):

**Assumption d.5.** The evolution of  $x(t)$  is given by the following system of continuous-time differential equations:

$$\dot{x}_i = D_i(x(t); s) = x_i(t) (1 - x_i(t)) \Phi_i(\cdot) + s_i (x_i(t)) \quad (2)$$

with  $\delta \geq 0$ ;  $\bar{\pi}_1 = \bar{\pi}_2 = \frac{1}{2}$  and  $\bar{\pi}_3 = \bar{\pi}_2$  (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  $\Phi$ ). Following Binmore and Samuelson (1996), this latter term is called drift: it opens the model to the possibility of a heterogeneity of behaviors. Binmore et al. (1995), derive an analogous system in the following way. At each point in time, a fixed proportion of players (of measure  $\frac{\delta}{1+\delta}$ ) is replaced by new individuals whose aggregate behavior is represented by a generic, constant, completely mixed strategy (i.e.  $\bar{\pi}_i$ ), while the rest of the population aggregate behavior follows the replicator dynamics. The relative importance of the drift is measured by  $\delta$ ; which we refer to as the drift level. We assume  $\delta$  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,  $\bar{\pi}_1$  and  $\bar{\pi}_2$  have been chosen to be 1/2 since only the value of  $\bar{\pi}_3$  turns out to be genuinely significant.

We start our analysis on the stability properties of NE looking at the case of the replicator dynamic without drift (i.e. when  $\delta = 0$ ): 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

The replicator dynamic and game  $\Gamma$

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  $\Phi$ : 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 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 the Appendix. 2

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

### Figure 3

The dynamic with drift and game  $\mu$

Let  $\mu \in (0; 1)$  be a generic element of the space of the feasible perturbations. Figure 3 shows the trajectories of the replicator dynamic with drift under two different specifications of  $\mu$ : Diagram 3b) represents a situation in which, in the proximity of  $NE^0$ , the drift against  $m_1^0$  is uniform across players, where in diagram 3a) the drift against  $m_3^0$  is lower. As the figures 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 figure 3b), while figure 3a) shows that 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.

We are interested in the convergence and stability properties of (2) when  $\mu \rightarrow 0$ ; considering two different configurations of the drift parameter  $\mu$ :

$$\begin{aligned} \text{CASE A : } \mu &= \frac{23i}{49} \frac{4^{p-30}}{30} \\ \text{CASE B : } \mu &= \frac{23i}{49} \frac{4^{p-30}}{30}; 1 \end{aligned}$$

Given  $\frac{23i}{49} \frac{4^{p-30}}{30} \approx 0.0222673$ , 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  $\hat{RE}(\mu)$  be the set of restpoints of (2) for  $\mu$  sufficiently close to 0. The following properties hold:

- a)  $\mu \in (0; 1)$ ;  $\hat{RE}(\mu)$  contains an element of  $NE^1$ ; which is also asymptotically stable.
- b) under CASE A  $\hat{RE}(\mu)$  contains also two additional restpoints, both belonging to  $NE^0$ ; one of which is asymptotically stable.

**Proof.** See the Appendix:2.

There is a striking similarity between the content of Proposition 6 and the findings of Binmore et al. (1995), as we pointed out in the introduction. They analyze the chain store game, in one of whose equilibrium components a player selects a weakly dominated strategy with positive probability. This component is interior attracting. Moreover, like our NE<sup>0</sup>; such component 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 claim. The intuition is similar to the one in Binmore et al. (1995). When there is drift, the strategies against which the weakly dominated strategy does poorly will have positive weight at all times and therefore the part of the dynamics that depend on payoffs pushes against the dominated strategy. But the 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 such dominated strategy). When the balance between these two forces is right, one gets a stable equilibrium with positive weight for the dominated strategy, as it happens in our example.

## 6 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. This result complements that obtained by Cabrales (1996) for the Abreu and Matsushima (1994) mechanism. Since Cabrales (1996) uses dynamics that are different from those used here, it would be interesting to check if the results we obtain here extend to Abreu and Matsushima (1994) games. More generally, a deeper study with evolutionary tools of other mechanisms studied in the literature would enhance our understanding of the performance of these mechanisms with boundedly rational agents, a necessary step before mechanisms are used in real life.

Ideally one would like to design a game for which convergence to the preferred social outcome could be guaranteed for the learning protocols that agents use. To achieve this goal, it is necessary to conduct empirical and experimental studies that reveal how people adjust their play in games like that studied in detail in this paper. The history of actual social arrangements may also give clues as to how people learn in such environments. Different mechanisms for public good provision have existed for centuries in many countries. These considerations imply the need for a substantial program of future research.

## References

- D. Abreu and H. Matsushima (1992), "Virtual Implementation in Iteratively Undominated Strategies: Complete Information", *Econometrica*, 60, 993-1008.



- D. Abreu and H. Matsushima (1994), "Exact Implementation", *Journal of Economic Theory*, 64, 1-19.
- K. Binmore, J. Gale and L. Samuelson (1995), "Learning to Be Imperfect: the Ultimatum Game", *Games and Economic Behavior*, 8, 56-90.
- K. Binmore and L. Samuelson (1996), "Evolutionary Drift and Equilibrium Selection", Institute for Advanced Studies, Vienna, Working Paper 26.
- T. Bärgers and R. Sarin (1993), "Learning through Reinforcement and Replicator Dynamics", UCL Discussion Paper 93-13.
- T. Bärgers (1994), "Weak Dominance and Approximate Common Knowledge", *Journal of Economic Theory*, 64, 265-276.
- A. Cabrales (1996), "Adaptive Dynamics and the Implementation Problem with Complete Information", Universitat Pompeu Fabra WP.
- E. Dekel and D. Fudenberg (1990), "Rational Behavior with Payo<sup>®</sup> Uncertainty", *Journal of Economic Theory*, 52, 243-267.
- T. Groves and J. Ledyard (1977), "Optimal Allocation of Public Goods: a Solution to the Free Rider Problem", *Econometrica*, 45, 783-809.
- M. O. Jackson (1992), "Implementation in Undominated Strategies: A Look at Bounded Mechanisms", *Review of Economic Studies*, 59, 757-775.
- M. O. Jackson, T. R. Palfrey and S. Srivastava (1994), "Undominated Nash Implementation in Bounded Mechanisms", *Games and Economic Behavior*, 6, 474-501.
- Y. G. Kim and J. Sobel (1995), "An Evolutionary Approach to Pre-Play Communication", *Econometrica*, 63, 1181-1193.
- E. Maskin (1977), "Nash Implementation and Welfare Optimality", mimeo, Massachusetts Institute of Technology.
- J. Moore (1990), "Implementation in Environments with Complete Information", in J. J. La<sup>®</sup>ont ed., *Advances in Economic Theory: Sixth World Congress*, Econometric Society.
- J. Moore and R. Repullo (1988), "Subgame Perfect Implementation", *Econometrica*, 58, 1083-1099.
- T. Muench and M. Walker (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.
- R. Repullo (1987), "A Simple Proof of Maskin's Theorem on Nash Implementation", *Social Choice and Welfare*, 4, 39-41.

- L. Samuelson (1993), "Does Evolution Eliminate Weakly Dominated Strategies?", in K. Binmore, A. Kirman and P. Tani (eds.) (1993) *Frontiers of Game Theory*, London, MIT Press.
- L. Samuelson and J. Zhang (1992), "Evolutionary Stability in Asymmetric Games", *Journal of Economic Theory*, 57, 363-391.
- K. Schlag (1994), "Why Imitate, and if so, How? Exploring a Model of Social Evolution", Friedrich-Wilhelms-Universität Bonn, Discussion Paper No. B-296.
- T. Sjöström (1994), "Implementation in Undominated Nash Equilibria without Integer Games", *Games and Economic Behavior*, 6, 502-511.
- PD Taylor and LB Jonker (1978), "Evolutionary Stable Strategies and Game Dynamics", *Mathematical Biosciences*, 40, 145-156
- P. de Trenchard (1988), "Stability of the Groves and Ledyard Mechanism", *Journal of Economic Theory*, 46, 164-171.
- J. Weibull (1995), *Evolutionary Game Theory*, Cambridge MA, MIT Press.

## 7 Appendix

**Proof of Proposition 1.** We already noticed that agent 3 has a weakly dominated strategy (namely,  $m_3^0$ ). In particular,  $m_3^1$  (truth-telling) makes agent 3 (strictly) better off than  $m_3^0$  (lying), unless agents 1 and 2 coordinate their actions completely, that is, unless they play  $m_i^0$   $i = 1, 2$  with probability 1 or they play  $m_i^1$   $i = 1, 2$  with probability 1, (in which case, 3 is completely indifferent). This leads to the following lemma:

**Lemma 1.** No strategy profile in which  $x_3 \in (0; 1)$  can be a Nash equilibrium unless  $x_1 = x_2 = 1$  or  $x_1 = x_2 = 0$ , that is, unless agents 1 and 2 play the same strategy with probability 1.

With this consideration in mind, we construct the proof as follows: we fix the mixed strategy of player 3 and check what mixed strategies for player 1 and 2 can sustain a Nash equilibrium. Noting that

$$U_{i=1}^1(m_1^1; x) - U_{i=1}^1(m_1^0; x) = \frac{1}{12}(x_2(x_3 - 1) + 7x_3 - 3) \quad (3)$$

$$U_{i=2}^1(m_2^1; x) - U_{i=2}^1(m_2^0; x) = \frac{1}{12}(x_1(x_3 - 1) + 7x_3 - 3) \quad (4)$$

we can make the following observations:

a) When  $x_3 < \frac{3}{7}$ ,  $m_i^0$  (lying) yields a strictly higher payoff than  $m_i^1$  for both 1 and 2, independently of what the other player does. Therefore, the strategy profiles in  $NE^0$  (and only those) will be Nash equilibria.

b) When  $x_3 = \frac{3}{7}$ ,  $m_i^0$  yields a strictly higher payoff than  $m_i^1$  unless  $x_2 = 0$ , and  $x_2 = 0$  makes player 1 indifferent between  $m_1^0$  and  $m_1^1$  (a symmetric argument holds for player 2). This excludes the possibility of  $(1; 1; \frac{3}{7})$  being a Nash equilibrium of the game, leaving  $(0; 0; \frac{3}{7}) \in NE^0$  as the unique Nash equilibrium when  $x_3 = \frac{3}{7}$ :

c) When  $x_3 \in (\frac{3}{7}; \frac{1}{2})$  there are no Nash equilibria. This happens because in this case if  $x_1 = 1$ , the best response of player 2 is  $x_2 = 0$  and if  $x_1 = 0$ , the best response for player 2 is  $x_2 = 1$ : However, neither  $(0; 1; x_3)$  nor  $(1; 0; x_3)$  can be Nash equilibria when  $x_3 \in (\frac{3}{7}; \frac{1}{2})$  by Lemma 1.

d)  $x_3 = \frac{1}{2}$ : In analogy with the case where  $x_3 = \frac{3}{7}$ , it is an implication of Lemma 1 that  $(1; 1; \frac{1}{2}) \in NE^1$  is the unique Nash equilibrium when  $x_3 = \frac{1}{2}$ :

e) When  $x_3 > \frac{1}{2}$  announcing  $m_i^1$  (truth-telling) is optimal for  $i = 1$  and 2, independently of what the other player does. Thus, the strategy profiles in  $NE^1$  (and only those) will be Nash equilibria.

Since this exhausts all cases the result follows.  $\square$

**Proof of Proposition 2.** To prove the proposition, it is enough to show that any interior trajectory converges. This is because, 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)) to obtain the result.

We start by observing that the fact that the dynamic is forward invariant 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 it must converge. This already implies convergence of player 3's mixed strategy. Let us denote  $x_i^\infty = \lim_{t \rightarrow \infty} x_i(t)$ , when such a limit exists. Three alternative cases have to be discussed:

- a)  $x_3^\infty = 0$ : If  $x_3^\infty = 0$  there must be a time  $t^0$  such that  $x_3(t) < \frac{3}{7}$  for  $t > t^0$ . This implies that there is a  $k > 0$  such that for all  $t^0 > t$ ,  $\Phi_{i,i}(x(t)) < k$  for  $i = 1; 2$ . This implies, by monotonicity,  $\lim_{t \rightarrow \infty} x_i(t) = 0$  for  $i = 1; 2$ , thus  $x^\infty = (0; 0; 0)$ .
- b)  $x_3^\infty = 1$ : By a similar argument, monotonicity implies  $x^\infty = (1; 1; 1)$ :
- c)  $x_3^\infty \in (0; 1)$ : We want to prove that  $x_3^\infty$  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  $x_3^\infty \in (0; 1)$  it then must be that both players 1 and 2 select, in the limit, the same pure strategy. Given that this result implies convergence of the entire mixed strategy profile, the result follows. More formally, what we have to prove is contained in the following lemma:

**Lemma 3.** If  $x_3^\infty \in (0; 1)$  then:

$$\begin{aligned} & \text{either} \\ & x_i^\infty = 0; i = 1; 2 \text{ (CASE 0 hereafter)} \\ & \text{or} \\ & x_i^\infty = 1; i = 1; 2: \text{ (CASE 1)} \end{aligned}$$

**Proof.** Let's assume, for the purpose of contradiction, that neither of the above statements is true. In that case, there must exist a sequence  $t_k, g_{k=1}^1$  and a positive constant  $\epsilon > 0$  such that either  $x_i(t_k) > \epsilon; i = 1; 2$  or  $x_i(t_k) < 1 - \epsilon; i = 1; 2$  for all  $k$  (in other words, the system must stay infinitely often an  $\epsilon$  away from the faces of  $\Phi$  in which player 1 and 2 play the same pure strategy). We already noticed that these are the only faces of  $\Phi$  in which both pure strategies for player 3 yield the same payoff: if the system stays away from them infinitely often along the solution path, it then must be that the cumulative payoff difference will grow unbounded as time goes to infinity. As we will see, this in turn implies (by monotonicity) that  $x_3(t)$  will also reach, in the limit, its highest value, that is,  $x_3^\infty = 1$ ; as a result of the extinction of the weakly dominated strategy  $m_3^0$ , which is a contradiction.

To show this, we first notice that the payoff difference  $\Phi_{i,i}(x(t))$  is a continuous function of  $x(t)$  defined over a compact set ( $\Phi$ ): In the case of player 3, such function takes the

following form:

$$\Phi_{13}(x(t)) = \frac{(x_1(t) - x_2(t))^2 + x_1(t)(1 - x_1(t)) + x_2(t)(1 - x_2(t))}{6} \quad (5)$$

Take  $g_M = \max_{i \in \{1,2\}; x_i \in \Phi_{13}} [g_i(m_i; x_{-i}(t))]$ , 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  $\Phi$ ): Then define  $\zeta_1; \zeta_2; \zeta_3$  and  $\zeta_4$  as follows:

$$\zeta_1 \text{ solves } 2 \exp[-g_M \zeta_1] = \frac{2}{3} \text{ (i.e. } \zeta_1 = \frac{\ln[2]}{g_M})$$

$$\zeta_2 \text{ solves } (1 - \frac{2}{3}) \exp[-g_M \zeta_2] = \frac{2}{3} \text{ (i.e. } \zeta_2 = \frac{\ln[1 - \frac{2}{3}]}{g_M})$$

$$\zeta_3 \text{ solves } 2 \exp[g_M \zeta_3] = 1 - \frac{2}{3} \text{ (i.e. } \zeta_3 = \frac{\ln[1 - \frac{1}{3}]}{g_M})$$

$$\zeta_4 \text{ solves } (1 - \frac{2}{3}) \exp[g_M \zeta_4] = 1 - \frac{2}{3} \text{ (i.e. } \zeta_4 = \frac{\ln[\frac{2}{3}]}{g_M})$$

and take  $\delta = \min\{\zeta_1; \zeta_2; \zeta_3; \zeta_4\}$ ; that is, set a lower bound for the time interval in which, after each  $t_k$ ,  $\frac{2}{3} < x_i < 1 - \frac{2}{3}$ ;  $i = 1; 2$  and therefore  $\Phi_{13}(x(t))$  still remains bounded away from 0 (i.e.  $\Phi_{13}(x(t)) > \frac{2(1 - \frac{2}{3})}{3} > 0$ ;  $\forall t \in [t_k; t_k + \delta]$ ): Denote by  $G_2 = \{x \in \Phi_{13}(x) \mid \frac{2(1 - \frac{2}{3})}{3} < \Phi_{13}(x) < \frac{2(1 - \frac{2}{3})}{3} + \delta\}$ . Now define:

$$\rho_i(x(t)) = \frac{\partial}{\partial t} \ln \frac{x_i(t)}{1 - x_i(t)} = \frac{x_i(t)}{x_i(t)} - \frac{(1 - x_i(t))}{1 - x_i(t)} = \frac{x_i(t)}{x_i(t) - (1 - x_i(t))^2}$$

i.e. 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. Notice that also  $\rho_3(x(t))$  will be a positive number bounded away from 0 infinitely often since, by assumption d.1, the difference in growth rates is a continuous function of  $x(t)$  defined on a compact set, which preserves the same sign of  $\Phi_{13}(x(t))$ : This implies that we can always define a constant  $g_2 = \min_{x \in G_2} \rho_3(x(t))$ , with  $g_2 > 0$  by assumption d.2. Also by assumption d.2  $\rho_3(x(t)) > g_2$  ( $\implies \Phi_{13}(x(t)) > \frac{2(1 - \frac{2}{3})}{3}$ ): If we integrate the value of  $\rho_3(x(t))$  over time we then obtain:

$$\lim_{t \rightarrow \infty} \int_0^t \rho_3(x(t)) dt = \sum_{k=1}^{\infty} \int_{t_k}^{t_k + \delta} \rho_3(x(t)) dt > g_2 \sum_{k=1}^{\infty} \int_{t_k}^{t_k + \delta} dt = \infty$$

which implies that  $x_3^\infty = 1$ ; which leads to a contradiction.  $\square$

To summarize, Lemma 3 shows that, if  $x_3^\infty \in (0; 1)$ ;  $x_1(t)$  and  $x_2(t)$  must converge (and therefore  $x(t)$  must converge to a Nash equilibrium). Since this exhausts all cases the result follows.  $\square$

**Proof of Proposition 3.**

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

$$u_i(\hat{x}_i; x_{-i}) = u_i(x_i; x_{-i}) \cdot \prod_{j \in i} x_j^{m_j} v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) + (1 - \prod_{j \in i} x_j^{m_j}) v_i(0; \hat{R}_i):$$

For all  $\hat{x}_i$ ,

$$u_i(\hat{x}_i; x_{-i}) \cdot \left(1 - \prod_{m_i \in S_i} \hat{x}_i\right) u_i(x_i; x_{-i}) + \prod_{m_i \in S_i} \hat{x}_i \cdot \left(\prod_{j \in i} x_j^{m_j} v_i(0; \hat{R}_i) + (1 - \prod_{j \in i} x_j^{m_j}) U_{in}\right):$$

Then

$$u_i(\hat{x}_i; x_{-i}) - u_i(x_i; x_{-i}) \cdot \prod_{m_i \in S_i} \hat{x}_i \cdot \left(\prod_{j \in i} x_j^{m_j} (v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) - v_i(0; \hat{R}_i)) + (1 - \prod_{j \in i} x_j^{m_j}) (v_i(0; \hat{R}_i) - U_{in})\right)$$

which is greater than zero since by the definition of  $k_j$ ,

$$\prod_{j \in i} x_j^{m_j} \geq \prod_{j \in i} k_j \geq \frac{U_{in} - v_i(0; \hat{R}_i)}{v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) - v_i(0; \hat{R}_i) + U_{in} - v_i(0; \hat{R}_i)}: \quad 2$$

**Proof of Proposition 4.** By contradiction.

Suppose that a) is the statement that stops being true earliest, that it does it for agent  $i$  and strategy  $m_i \in S_i$  and that the boundary time is  $t^0$ . Then it must be true that

$$\frac{x_i^{m_i}(t^0)}{x_i^{m_i}(0)} = \exp[h_i t^0] \frac{H}{x_i^{m_i}(0)}$$

Notice that for all  $t$

$$\begin{aligned} u_i(x_i^{m_i}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) \cdot \left(\prod_{j \in i} x_j^{m_j}(t) v_i(0; \hat{R}_i) + U_i (1 - \prod_{j \in i} x_j^{m_j}(t))\right) \\ - \left(\prod_{j \in i} x_j^{m_j}(t) (v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) - v_i(0; \hat{R}_i)) + (1 - \prod_{j \in i} x_j^{m_j}(t)) (v_i(0; \hat{R}_i) - U_{in})\right) \\ = [U_i - v_i(0; \hat{R}_i) - \prod_{j \in i} x_j^{m_j}(t) (v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) + U_i - 2v_i(0; \hat{R}_i))] \end{aligned}$$

But since b) is true for  $t < t^0$

$$u_i(x_i^{m_i}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) < \prod_{j \in i} x_j^{m_j}(t) \cdot \left( H^{n_i-1} (v_i(f_i(\hat{A}(i; R^a)); \hat{R}_i) + U_i - 2v_i(0; \hat{R}_i)) \right)$$

So we have that

$$u_i(x_i^{m_i}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) < h_v$$

which by assumption d.2 and the definition of  $h_v$  and  $h_g$  implies that

$$g_i(m_i; x_{-i}(t)) - g_i(m_i^a; x_{-i}(t)) < h_g:$$

which integrating from 0 to  $t^0$  and given that  $x_i^{m_i^a}(t^0) \cdot H$  implies that

$$\frac{x_i^{m_i}(t^0)}{x_i^{m_i}(0)} < \exp[h_g t^0] \frac{H}{x_i^{m_i}(0)}$$

This is a contradiction.

Suppose that b) is the statement that stops being true earliest, that it does it for agent  $i$  and that the boundary time is  $t^0$ . Then it must be true that  $x_i^{m_i^a}(t^0) = H$ .

First notice that for all  $m_i \in S_i$  since the payoffs of strategy  $m_i^a$  and other strategies in  $S_i$  differ only when playing against strategies not in  $S_i$

$$u_i(x_i^{m_i^a}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) = U_i \prod_{j \in I} x_j^{s_j}(t)$$

since a) holds for  $t < t^0$

$$u_i(x_i^{m_i^a}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) > U_i \exp[h_g t] \frac{H}{x_i^{m_i}(0)} \prod_{j \in I} x_j^{s_j}(0)$$

Since  $U_i \frac{H}{x_i^{m_i}(0)} \prod_{j \in I} x_j^{s_j}(0) < \mu$ , this implies by assumption d.3 that

$$(g_i(m_i^a; x_{-i}(t)) - g_i(m_i; x_{-i}(t))) > \ln \left[ U_i \frac{H}{x_i^{m_i}(0)} \prod_{j \in I} x_j^{s_j}(0) \right] + h_g t$$

So by integration we have that

$$\frac{x_i^{m_i^a}(t^0)}{x_i^{m_i^a}(0)} \frac{x_i^{m_i}(0)}{x_i^{m_i}(t^0)} > \exp \left[ \frac{1}{h_g} \left( \ln \left[ U_i \frac{H}{x_i^{m_i}(0)} \prod_{j \in I} x_j^{s_j}(0) \right] + h_g t^0 \right) \right] = L$$

Adding over all strategies in  $S_i$  we have

$$\frac{x_i^{m_i^a}(t^0)}{x_i^{m_i^a}(0)} > \frac{x_i^{S_i}(t^0)}{x_i^{S_i}(0)} L = \frac{1}{1} \frac{x_i^{S_i}(t^0)}{x_i^{S_i}(0)} L > L$$

But this implies that  $x_i^{m_i^a}(t^0) > H$  (using the assumption that  $x_i^{m_i^a}(0) L > H$ ), which is a contradiction.

Suppose that c) is the statement that stops being true earliest, that it does it for agent  $i$  and that the boundary time is  $t^0$ . Then it must be true that  $\frac{x_i^{m_i^a}(t^0)}{x_i^{m_i^a}(0)} = L \frac{x_i^{m_i^a}(0)}{x_i^{m_i^a}(0)}$

As before, notice that for all  $m_i \in S_i$  the payoffs of strategy  $m_i$  and  $m_i$  differ only when playing against strategies not in  $S_i$ ; so

$$u_i(x_i^{m_i}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) \cdot U_i(x_i^{m_i}(t); x_{-i}(t))$$

which by part a) of the proposition implies that for all  $t < t^0$

$$u_i(x_i^{m_i}(t); x_{-i}(t)) - u_i(x_i^{m_i}(t); x_{-i}(t)) \cdot U_i(x_i^{m_i}(0); x_{-i}(0)) \exp\left[-\int_0^t h_j x_j^{s_j} dt\right]$$

Since  $U_i(x_i^{m_i}(0); x_{-i}(0)) < \mu$ , this implies by assumption d.3 that

$$(g_i(m_i; x_{-i}(t)) - g_i(m_i; x_{-i}(t))) > \ln U_i(x_i^{m_i}(0); x_{-i}(0)) + h_g t$$

So by integration we have that

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

which is a contradiction. Since this exhausts all cases the result follows.  $\square$

### Proof of Proposition 5.

(i) We know, from Proposition 2, that  $x_3 > 0$  in any interior point. This implies that if there is a time  $t$  such that  $x_3(t) > \frac{1}{2}$ , then  $x_3(t^0) > \frac{1}{2}$  for all  $t^0 > t$ . From equations (3-4) we have that, whenever  $x_3(t) > \frac{1}{2}$ ,  $\Phi_i(x) > 0$  for players 1 and 2. This implies that if there is a time  $t$  such that  $x_3(t) > \frac{1}{2}$ , then  $x_i > 0$  for all  $t^0 > t$  for  $i = 1, 2$  and therefore  $x(t)$  converges. Since convergence must be to a Nash equilibrium and  $x_1$  and  $x_2$  have been increasing,  $x$  converges to  $NE^1$ . To show the stability of  $NE^1$  it suffices to show that there is a neighborhood of  $NE^1$  such that, for all  $x(0)$  in this neighborhood, there is a time  $t$  such that  $x_3(t) > \frac{1}{2}$ . Let  $x_i(0) = 1 - \alpha_i$  for  $i = 1, 2$  and  $x_3(0) = \frac{1}{2} - \beta$ , with  $\alpha_i > 0, \beta > 0$ . From (3-4) we also have that  $1 - \alpha_i < \Phi_i(x) < 1$  for  $i = 1, 2$ , thus

$$\exp[\alpha_i t](1 - \alpha_i) < x_i(t) < \exp[t](1 - \alpha_i); \tag{6}$$

Since  $\Phi_3(x) = \frac{x_1(1 - x_1)}{6}$  we have by equation (6)

$$\frac{x_3(t)}{x_3(0)} > \frac{(1 - \alpha_1)(\exp[\alpha_1 t](1 - \alpha_1))}{6}$$

thus

$$\frac{x_3(t)}{x_3(0)} > \frac{(1 - \alpha_1)(\exp[\alpha_1 t] - (1 - \alpha_1))}{6} > \frac{(1 - \alpha_1)(\alpha_1 t + \alpha_1)}{6}$$



This implies that

$$x_3(t) > \exp \frac{2}{4} \frac{(1 - \alpha_1)(\frac{t^2}{2} + \alpha_1 t)}{6} \frac{3}{5} \left(\frac{1}{2} - \alpha_1\right):$$

Note that for  $t = \alpha_1$

$$\exp \frac{2}{4} \frac{(1 - \alpha_1)(\frac{t^2}{2} + \alpha_1 t)}{6} \frac{3}{5} \exp \frac{2}{4} \frac{(1 - \alpha_1)(\frac{\alpha_1^2}{2})}{6} \frac{3}{5} > 1$$

and therefore  $x_3(t) > \frac{1}{2}$  for  $\alpha_1$  small enough, which is what we wanted to show.

(ii). Assume that  $x_3(0) > \frac{3}{7}$ . Since  $\dot{x}_3(t) \geq 0$  for all  $t$ ,  $x_3(t)$  is an increasing function of  $t$ , therefore it must converge. Since the initial condition  $x_3(0)$  is larger than  $\frac{3}{7}$  it must converge to a number larger than  $\frac{3}{7}$ . We know that  $x(t)$  converges to a Nash equilibrium by Proposition 2. Since there is no equilibrium in  $NE^0$  with  $x_3 > \frac{3}{7}$ ,  $x(t)$  cannot converge to a point in  $NE^0$ . Since  $x_3(0)$  can be arbitrarily close to  $\frac{3}{7}$  and therefore to the set  $NE^0$ , this set must be unstable.  $\square$

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

We start by observing that, given  $\epsilon \in (0, 1)$ ; any rest point must be completely mixed, and it also must be  $x_3 > \epsilon$ ; as  $\Phi_{i3}(\cdot)$  is always positive in the interior of the state space  $\Phi$  (because  $m_3^0$  is a weakly dominated strategy). We also know, by the continuity of the vector field with respect to  $\epsilon$ ; that every limiting rest point of the dynamic, as  $\epsilon$  goes to zero, must lie in the set of restpoints of the unperturbed dynamic  $RE(\epsilon)$ .

We analyze first the limit set of rest points under CASE 0. In this case, both players 1 and 2 play their strategy  $m_i^0$  with probability 1, that is  $x_i^0 = 0$ ; for  $i = 1, 2$ . Setting  $\underline{x}_1 = 0$  yields the following equation:

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

and an analogous expression can be obtained for  $\underline{x}_2$ : Denote by  $x_3^0$  a limiting value in a rest point, if a limit exists, for  $x_3$ . When the limiting values for  $x_1$  and  $x_2$  are zero we have:

$$\lim_{\epsilon \rightarrow 0} \frac{x_i}{\epsilon} = \frac{1}{2(3 - 7x_3^0)} \quad (8)$$

Notice that in this case if a rest point exists it must be  $x_3^0 < \frac{3}{7}$ ; since  $\underline{x}_1 > 0$ : We then set  $\underline{x}_3 = 0$ ; substitute  $\underline{x}_i$  with the expression in (8), solve for  $x_3$ ; and substitute  $x_i$ ;  $i = 1, 2$  and  $\epsilon$  by their limiting value of zero. The solutions for  $x_3^0$  take the following form:

$$x_3^0 = \frac{1 + 7\epsilon + \sqrt{1 - (46\epsilon - 49\epsilon^2)}}{10} \quad \text{and} \quad x_3^0 = \frac{1 + 7\epsilon - \sqrt{1 - (46\epsilon - 49\epsilon^2)}}{10}$$

Remember that  $x_3^0$  must be a real, positive number, with  $0 < x_3^0 < \frac{3}{7}$ . For the expression under the square root at the numerator to be nonnegative, it must be that  $\frac{23i-4}{49} \geq \frac{1}{4} : 0.222673$ ; which determines the feasible range for both roots. Within this interval of values for  $\frac{23i-4}{49}$ ,  $x_3^0$  ( $x_3^0$ ) is a strictly decreasing (increasing) function of  $\frac{23i-4}{49}$ ; which has a minimum and a maximum, whose values are  $\frac{15i-2}{35}$  (0) and  $\frac{2}{10} - \frac{15i-2}{35}$  respectively. As  $\frac{23i-4}{49} \rightarrow \frac{1}{4}$ ; both solutions converge to  $\frac{15i-2}{35}$ .

We now deal with the subset of limiting rest points under CASE 1, i.e. with limiting values for  $x_i = 1$  for  $i = 1; 2$ . The equations corresponding to (7-8) are now the following:

$$\frac{(1 - x_1)}{x_1} = \frac{\frac{1}{2} x_1}{x_1 \frac{1}{3} + (1 - x_3) \left( \frac{1}{12} (1 - x_2) - \frac{2}{3} \right)} \quad (9)$$

$$\lim_{x_i \rightarrow 1} \frac{(1 - x_1)}{x_1} = \frac{1}{2 - \frac{1}{3} - \frac{2}{3} (1 - x_3)} \quad (10)$$

Denote by  $x_3^1$  a limiting value in a rest point for  $x_3$  in this latter case: By analogy with CASE 0, we know from (10) that, if a rest point exists, it must be  $x_3^1 > \frac{1}{2}$ : There is a unique feasible solution for  $x_3^1 \in (\frac{1}{2}, 1)$  with the following form:

$$x_3^1 = \frac{3 + 4\sqrt{9 - 16(1 - \frac{1}{2})}}{10}$$

Following the same procedure for the remaining rest points of the unperturbed dynamics (i.e. the pure strategy profiles which belong to RE ( $\mu$ ) and do not satisfy either CASE 0 or CASE 1) does not add any element to the limiting set of rest points of the perturbed dynamics. This should not be surprising, as any other rest point 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 for the dynamic system is as follows:

$$J(x; \mu) = \begin{pmatrix} (1 - 2x_1)\Phi_{11} & \frac{i(1-x_1)x_1(1-x_3)}{12} & \frac{(1-x_1)x_1(7+x_2)}{12} \\ \frac{i(1-x_2)x_2(1-x_3)}{12} & (1 - 2x_2)\Phi_{22} & \frac{(1-x_2)x_2(7+x_1)}{12} \\ \frac{(1-2x_2)(1-x_3)x_3}{6} & \frac{(1-2x_1)(1-x_3)x_3}{6} & (1 - 2x_3)\Phi_{33} \end{pmatrix}$$

We analyze CASE 0 first. We know that, in this case, we have two restpoints, which we call  $x^0 = (0; 0; x_3^0)$  and  $x^0 = (0; 0; x_3^0)$ : We evaluate the Jacobian when  $x_1, x_2$  and  $\mu$  are equal to their limiting value (i.e. 0). The corresponding eigenvalues are:  $0; \frac{i(3+7x_3^0)}{12}; \frac{i(3+7x_3^0)}{12}$ : 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 the continuity of  $J(\cdot)$  ensures that the other two will be negative, for any  $\mu$ .

sufficiently small. We now linearize the rest points (as a function of  $\mu$ ) around  $NE^0$ . We set  $x(\mu; \pm) = (\pm_1 \mu; \pm_2 \mu; x_3^0 + \pm_3 \mu)$ ; where  $\pm = (\pm_1; \pm_2; \pm_3)$  denotes the vector collecting the coefficients of the linearised system. We then evaluate the following expression:

$$\hat{A}^0(x_3^0; \pm) = \lim_{\mu \rightarrow 0} \frac{\det J(x; \mu)_{x(x; \pm)}}{\mu^3}$$

We do so because  $\det(J(x; \mu))$ ; which is equal to zero  $8 \times 2$   $NE^0$ ; 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  $\mu$  is sufficiently small. For CASE 0 we get the following result:

$$\hat{A}^0(x_3^0; \pm) = \frac{\mu^3 [54 + x_3^0(252 + 294x_3^0) + (\pm_1 + \pm_2) 9 \mu + 39x_3^0 + 63(x_3^0)^2 + 49(x_3^0)^3]}{864} \quad (11)$$

We first notice that (11) does not depend on  $\pm_3$ . To evaluate  $\text{sign}(\hat{A}^0(x_3^0; \pm))$  we only need to get estimates of  $\pm_1$  and  $\pm_2$ , the linear coefficients which measure the responsiveness of the equilibrium values of  $x_i$ ;  $i = 1; 2$  to small changes in  $\mu$ . We do so setting  $\lim_{\mu \rightarrow 0} \frac{\partial}{\partial \mu} D(x; \mu)_{x(x; \pm)} = 0$  and solving for  $f_{\pm_1; \pm_2; x_3^0}$ : There are two alternative set of solutions, each of them corresponds to each of the restpoints. In particular:

$$\begin{aligned} \hat{\pm}_1^0 = \hat{\pm}_2^0 &= \frac{23 \mu + 49 \mu + 7 \mu \sqrt{1 \mu - (46 \mu + 49 \mu)}}{8} \\ \hat{\pm}_1^0 = \hat{\pm}_2^0 &= \frac{23 \mu + 49 \mu + 7 \mu \sqrt{1 \mu - (46 \mu + 49 \mu)}}{8} \end{aligned}$$

We evaluate the numerator of (11) for both sets of solutions, and we get the following expressions:

$$\hat{A}^0(-) = \frac{\mu^3 [7 + 322 \mu + 343 \mu^2 + (49 \mu + 23) \sqrt{1 \mu + 46 \mu + 49 \mu^2}]}{10} \quad (12)$$

$$\hat{A}^0(+) = \frac{2863 \mu + 147476 \mu^2 + 882882 \mu^3 + 1546244 \mu^4 + 823543 \mu^5 + k \sqrt{1 \mu + 46 \mu + 49 \mu^2}}{1000} \quad (13)$$

with  $k = (3887 \mu + 60123 \mu^2 + 165669 \mu^3 + 117649 \mu^4)$ :

Both  $\hat{A}^0(-)$  and  $\hat{A}^0(+)$  are plotted in Figure 4. As the diagram shows,  $\hat{A}^0(-)$  is always negative in the domain  $0; \frac{23 + \sqrt{49 - 30}}{49}$ ; whereas  $\hat{A}^0(+)$  is not. As a result of that,  $x^0$  is asymptotically stable whereas  $\hat{x}^0$  is not.

Figure 4

Asymptotic stability of the dynamic with drift

We now move on to CASE 1. Here we have a unique rest point, which we call  $\hat{x}^1$  (1; 1;  $\hat{x}_3^1$ ). The eigenvalues of the unperturbed dynamics are as follows:  $0; \frac{1-2x_3}{3}; \frac{1-2x_3}{3}$ . As in CASE 0, there are two (identical) negative eigenvalues (given that  $x_3 > \frac{1}{2}$ ); and the remaining eigenvalue equal to zero. By analogy with CASE 0, we now define  $x(\epsilon) = (1 + \epsilon_1; 1 + \epsilon_2; x_3^0 + \epsilon_3)$  and solve  $\lim_{\epsilon \rightarrow 0} \frac{D(x; \epsilon)}{x(\epsilon; \epsilon)} = 0$  to get estimates of  $\epsilon$ . The unique feasible solution (corresponding to the unique limiting equilibrium), takes the following form:

$$\hat{\epsilon}_1 = \hat{\epsilon}_2 = \frac{3 \cdot 2 \cdot 4^{-3} + \rho \cdot 9 \cdot 16^{-1} + 16^{-2}}{2}$$

The function corresponding to (12-13) takes now the following form:

$$\hat{A}^1(\epsilon) = \frac{24 \cdot (1 + \epsilon)^{\rho} + (2 \cdot 4^{-3}) \cdot \rho \cdot \epsilon^{\rho}}{5}$$

with  $\rho = 9 \cdot 16^{-1}$ . The function  $\hat{A}^1(\epsilon)$  is also plotted in Figure 4. As the diagram shows,  $\hat{A}^1(\epsilon)$  stays negative  $\forall \epsilon \in (0; 1)$ . As a result of that,  $\hat{x}^1$  is asymptotically stable under any drift configuration. 2

$m_3 = m_3^0$	$m_3 = m_3^1$												
$m_2^0$ $m_2^1$	$m_2^0$ $m_2^1$												
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 2px;"><math>m_1^0</math></td> <td style="padding: 2px;"><math>\frac{1}{4}, \frac{1}{4}, \frac{1}{2}</math></td> <td style="padding: 2px;"><math>\frac{1}{3}, 0, \frac{1}{3}</math></td> </tr> <tr> <td style="padding: 2px;"><math>m_1^1</math></td> <td style="padding: 2px;"><math>0, \frac{1}{3}, \frac{1}{3}</math></td> <td style="padding: 2px;"><math>0, 0, \frac{1}{3}</math></td> </tr> </table>	$m_1^0$	$\frac{1}{4}, \frac{1}{4}, \frac{1}{2}$	$\frac{1}{3}, 0, \frac{1}{3}$	$m_1^1$	$0, \frac{1}{3}, \frac{1}{3}$	$0, 0, \frac{1}{3}$	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 2px;"><math>m_1^0</math></td> <td style="padding: 2px;"><math>0, 0, \frac{1}{2}</math></td> <td style="padding: 2px;"><math>0, \frac{1}{3}, \frac{1}{2}</math></td> </tr> <tr> <td style="padding: 2px;"><math>m_1^1</math></td> <td style="padding: 2px;"><math>\frac{1}{3}, 0, \frac{1}{2}</math></td> <td style="padding: 2px;"><math>\frac{1}{3}, \frac{1}{3}, \frac{1}{3}</math></td> </tr> </table>	$m_1^0$	$0, 0, \frac{1}{2}$	$0, \frac{1}{3}, \frac{1}{2}$	$m_1^1$	$\frac{1}{3}, 0, \frac{1}{2}$	$\frac{1}{3}, \frac{1}{3}, \frac{1}{3}$
$m_1^0$	$\frac{1}{4}, \frac{1}{4}, \frac{1}{2}$	$\frac{1}{3}, 0, \frac{1}{3}$											
$m_1^1$	$0, \frac{1}{3}, \frac{1}{3}$	$0, 0, \frac{1}{3}$											
$m_1^0$	$0, 0, \frac{1}{2}$	$0, \frac{1}{3}, \frac{1}{2}$											
$m_1^1$	$\frac{1}{3}, 0, \frac{1}{2}$	$\frac{1}{3}, \frac{1}{3}, \frac{1}{3}$											

Figure 1

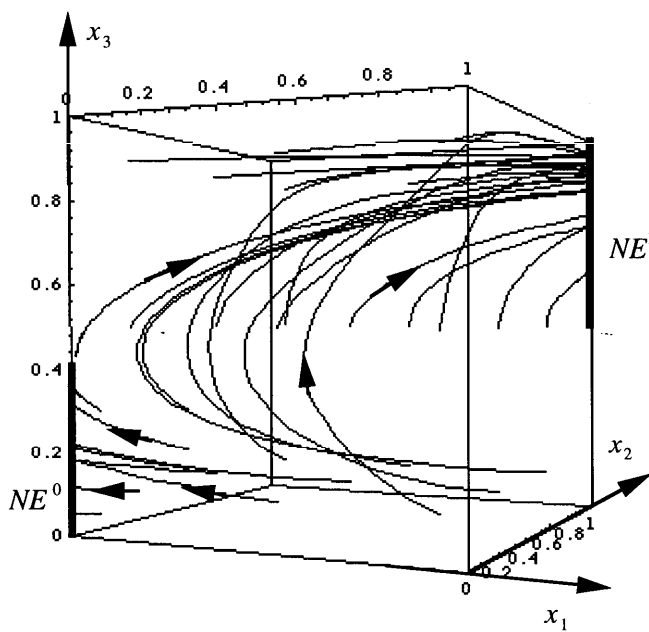


Figure 2

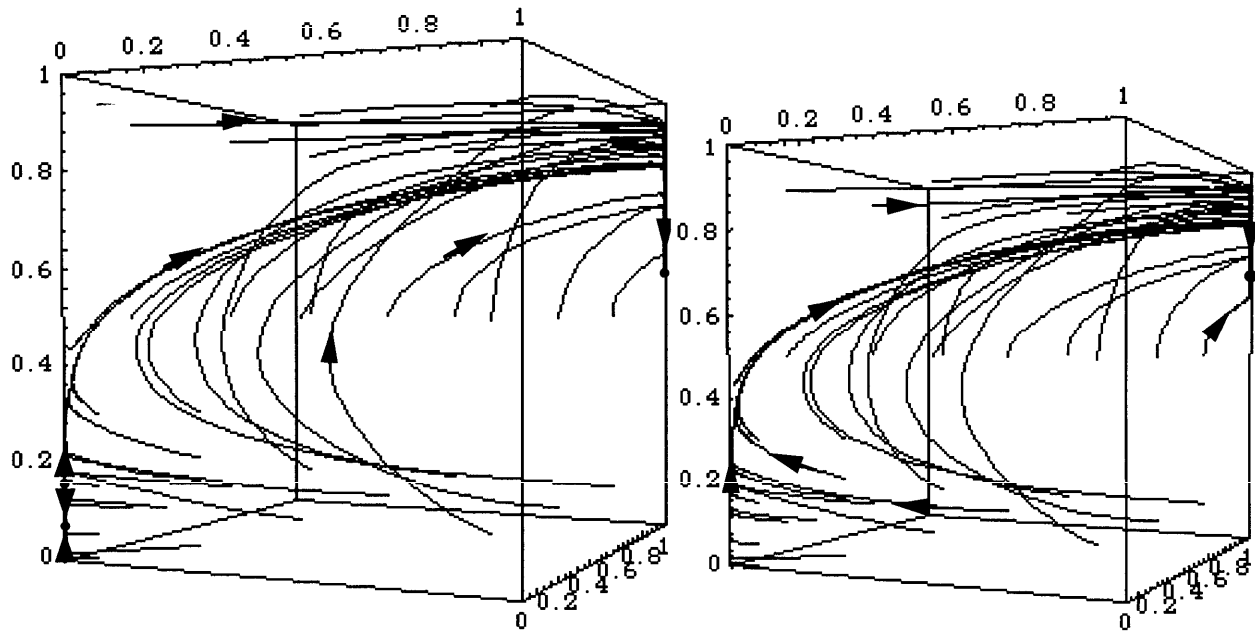


Figure 3

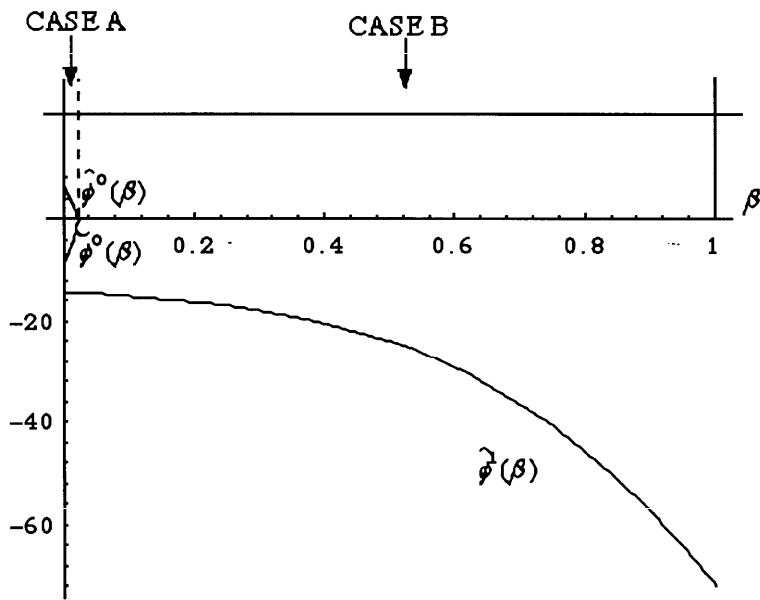


Figure 4