

## EVOLUTIONARY SELECTION WITH DISCRIMINATING PLAYERS

Abhijit BANERJEE<sup>1\*</sup>  
Jörgen W. WEIBULL<sup>2\*\*</sup>

Forschungsbericht/  
Research Memorandum No. 318

March 1993

First draft June 1992. Current version 17 February 1993.

\* Department of Economics, Harvard University, Cambridge MA 02138, USA

\*\* Department of Economics and Institute for International Economic Studies, Stockholm University,  
10 691 Stockholm, Sweden

<sup>1</sup> Both authors thank David Cooper, Drew Fudenberg, Vijay Krishna and Asher Wolinsky for helpful comments. Banerjee thanks the Institute for International Economic Studies, Stockholm University, for its hospitality during part of this research project.

<sup>2</sup> The research of Weibull was sponsored by the Industrial Institute for Economic and Social Research, Stockholm, Sweden. Weibull thanks the Institute for Advanced Studies, Vienna, for its hospitality during part of this research project.

Die in diesem Forschungsbericht getroffenen Aussagen liegen im Verantwortungsbereich des Autors/der Autorin (der Autoren/Autorinnen) und sollen daher nicht als Aussagen des Instituts für Höhere Studien wiedergegeben werden. Nachdruck nur auszugsweise und mit genauer Quellenangabe gestattet.

-----

All contributions are to be regarded as preliminary and should not be quoted without consent of the respective author(s). All contributions are personal and any opinions expressed should never be regarded as opinion of the Institute for Advanced Studies.

This series contains investigations by the members of the Institute's staff, visiting professors, and others working in collaboration with our departments.

## **Abstract**

This paper studies evolutionary games in which players can condition their strategy choice on some observable characteristic of their opponent, a characteristic we call their type. Recently, examples have been provided in which some players discriminate in this way, causing the evolutionary process to converge on non-Nash equilibrium play. Moreover, in some cases this generalization of the standard set-up of evolutionary game theory has been shown to destabilize certain inefficient Nash equilibria. We here provide a general model of evolutionary selection among discriminating behaviors, and find that the above examples are not robust; the close connection between evolutionary selection and Nash equilibrium, already established for the standard set-up, continues to hold, albeit in a slightly more complex form. Moreover, inefficient Nash equilibria may indeed be (weakly) stable in the evolutionary dynamics, and efficient Nash equilibria may be unstable.

## **Zusammenfassung**

Diese Studie befaßt sich mit evolutionären Spielen, in welchen die Spieler ihre Wahl von Strategien auf eine beobachtbare Charakteristik ihrer Opponenten bedingen können, eine Charakteristik, die wir ihren *Typus* nennen. In den letzten Jahren sind Beispiele von Spielen vorgestellt worden, in welchen einige Spieler in diesem Sinne diskriminieren können, mit dem Resultat, daß die evolutionären Prozesse zu einem Zustand tendieren, der kein Nash Gleichgewicht ist. In einigen Spielen hat man gezeigt, daß diese Verallgemeinerung der konventionellen Formulierung der evolutionären Spieltheorie einen destabilisierenden Effekt auf einige nicht effiziente Nash Gleichgewichte hat. In dieser Studie wird ein allgemeines Modell evolutionärer Auswahl zwischen verschiedenen diskriminierenden "behaviors" entwickelt, und wir zeigen, daß die genannten Beispiele nicht robust sind; das enge Verhältnis zwischen evolutionärer Auswahl und Nash Gleichgewicht, bekannt von früheren Studien auf diesem Gebiet, bleibt bestehen, nur in einer etwas komplizierteren Form. Übrigens: ineffiziente Nash Gleichgewichte können (Lyapunov) stabil sein in solchen evolutionären Auswahlprozessen, und effiziente Nash Gleichgewichte können instabil sein.



# 1 Introduction

This paper studies symmetric two-person evolutionary games with discriminating players. In contrast with standard evolutionary games, where each player always plays the same strategy, we introduce a framework in which players can change their strategies in response to some observable characteristic of their opponent, a characteristic we call their *type*.

A number of recent papers have already pointed out some of the more striking possibilities deriving from this kind of discriminating behavior. For example, Robson [14] and Banerjee and Weibull [1] show that the evolutionary process need not lead to Nash equilibrium outcomes when players condition their strategy choice on the type of opponent. This contrasts with the general result for the standard setting of evolutionary game theory (where each player uses the same strategy against all opponents) that the evolutionary processes leads the population state towards Nash equilibrium play, see *e.g.* Banerjee and Weibull [2] for a recent survey. Moreover, a number of recent studies of pre-play communication in co-ordination games show that the evolutionary process may lead the population state away from inefficient Nash equilibria when players condition their actions on pre-play messages, see *e.g.* Wärneryd [19] and Kim and Sobel [10].

The first set of examples raises some potentially serious questions about the evolutionary justification of Nash equilibrium. After all, there is nothing *per se* implausible about discriminating behavior, and a robust justification of Nash equilibrium should be able to accommodate this and other reasonable variations of the original evolutionary model. The second set of examples are more ambiguous; they have been interpreted as saying that the evolutionary mechanism selects *among* Nash equilibria, in favor of efficient or risk dominant equilibria. However, it is equally plausible to draw the implication that the evolutionary process selects against inefficient *outcomes* irrespective of whether these are Nash equilibria or not.

The aim of the present paper is to try to understand the general principles behind these examples. To this end we formulate a general model of evolutionary selection among discriminating behaviors in symmetric two-person games and study its general dynamic properties. The standard assumption in evolutionary game theory that each individual uses one and the same strategy against all opponents becomes a special case of this setting.

Our analysis of this generalized model reveals that the close connec-

tion between evolutionary selection and Nash equilibrium play, suitably re-interpreted, continues to hold when players are allowed to condition their strategy choice on the type of their opponent. In particular, we find that any *Lyapunov* (or *weakly stable*) stationary outcome of the evolutionary selection process, as operating on the full simplex of all discriminating behaviors, involves only Nash equilibrium play. More precisely, in such states the pure-strategy distributions in interactions within each player type, as well as between any two player types, always constitute Nash equilibria. Hence, in this generalized setting, all Lyapunov stable states are convex combinations of Nash equilibria of the underlying game.

This general result immediately tells us that the non-Nash outcomes derived in the examples above are not stable when *all* discriminatory behaviors are allowed for. Even though these outcomes are stable when the dynamics is restricted to some boundary face of the full simplex of all discriminating behaviors for the game, as in the mentioned examples, they are unstable in the unrestricted dynamics on the full simplex.

Moreover, we find that the set of stationary payoffs expands as more types are introduced, and that this set converges to a limit set  $V_\infty$  which is dense in the set  $U^{NE}$  of "symmetric convex combinations" of Nash equilibrium payoffs of the underlying game (see Section 5 for definitions). Also the subset of Lyapunov (or weakly) stable stationary payoffs converges towards a limit set as more types are introduced. However, this limit set,  $W_\infty \subset V_\infty$  may be a proper subset of the set  $U^{NE}$  in the strict sense of differing by at least some interval of positive length. (In particular,  $W_\infty$  need not be dense in  $U^{NE}$ .) In standard  $2 \times 2$  co-ordination games the limiting stable payoff set  $W_\infty$  consists of the "good" *and* the "bad" (strict) Nash equilibrium payoffs, as well as some, but not all, convex combinations of these two numbers. In particular, there is a whole interval of convex combinations near the "bad" equilibrium payoff which contains no stable payoff.

The fact that the limiting payoff set contains the "bad" Nash payoff tells us that, at least in the present formulation, there is no evolutionary tendency leading away from a (strict) Nash equilibrium which is Pareto dominated by another Nash equilibrium. However, the lack of such a "drift" away from inefficiency is not surprising, given what is currently known about evolutionary stability in co-ordination games with pre-play communication (see *e.g.* Bhaskar [4]). In the present set-up, "types" correspond to "messages" in these communication games. For such communication to destabilize an

inefficient outcome of the underlying game it is crucial that there exists a message (type) which is unmet in (absent from) the stationary state — this message is used to destabilize the state. To generate a (Lyapunov) stable but inefficient outcome in our framework, all we essentially need to do is to have all types present in the corresponding stationary state and let all players play the inefficient (strict Nash) equilibrium strategy with each other.

The basic result - that Lyapunov stable states are convex combinations of Nash equilibria - is also valid in more general settings. In particular, it continues to hold if types carry (fixed) costs and it is costly to distinguish types. However, when such costs are introduced, the inefficient (strict) Nash equilibrium in co-ordination games is no longer stable, since the more costly types are selected against in the evolutionary process, and if there is an absent (costly) player type in an inefficient stationary state, then the potential entry of this type destabilizes the stationary state (granted these costs are small). Therefore, there cannot be a stable stationary state in which the average payoff coincides with the “bad” Nash equilibrium of the underlying co-ordination game. Moreover, it turns out that the same argument, based on costly types and costly type-identification, rules out the existence of *any* (Lyapunov) stable state in games of the Prisoner’s Dilemma variety. In fact, any game which has a Nash equilibrium which is unique, symmetric, pure, and Pareto dominated by some other symmetric strategy combination has the same non-existence property. Thus, there are many games which have inefficient stable states in the standard evolutionary model but which lack stable states in the present generalized model with discriminating behaviors. It thus appears that discriminating behavior introduces a tendency away from inefficient outcomes, and, in certain games, thereby destroying the evolutionary tendency towards Nash equilibrium behavior.

The plan of the paper is as follows. In Section 2 we present the motivating examples. The formal model is developed in Section 3, and the basic results are derived in Section 4. Section 5 describes the limiting properties of the sets of stationary and stable payoffs, respectively, and Section 6 concludes with a summary and a brief discussion of the case in which different types carry different costs and the capacity to distinguish types is costly.

## 2 Examples

Figure 1 (a) illustrates a standard Prisoner's Dilemma game. In the replicator dynamics (in fact, in any monotonic dynamics), as applied to the standard set-up of evolutionary games, the population state converges from any interior initial state to the state in which all individuals in the population play the dominant strategy 2 ("defect"). Consider, however, a situation with discriminating behavior, as suggested in Robson [14], in which there are three types of player, labeled 1, 2 and 3. Individuals of type 1 always play strategy 1 ("cooperate") and type-2 individuals always play "defect", while individuals of type 3 play "defect" against types 1 and 2, but "cooperate" against other individuals of type 3. The payoffs to these three types of player are given in Figure 1 (b).

$$\begin{pmatrix} 1 & -1 \\ 2 & 0 \end{pmatrix} \quad \begin{pmatrix} 1 & -1 & -1 \\ 2 & 0 & 0 \\ 2 & 0 & 1 \end{pmatrix}$$

Figure 1: A Prisoner's Dilemma game.

Notice that once we have constructed the payoff matrix in Figure 1 (b) we can treat it as the payoff matrix of a standard evolutionary game with non-discriminating players who happen to have these payoffs. Therefore, all results in standard evolutionary game theory apply to this partially expanded game. In particular, the strictly dominated strategy 1 will be eliminated in the long run. Strategy 3 will not decrease over time since it is (weakly) dominant, so also strategy 2, which is weakly dominated by strategy 3, will be eliminated in the evolutionary process. Hence, only strategy 3 survives in the long run, implying that evolutionary selection will result in all players playing "cooperate", which is not a Nash equilibrium strategy in the original game (see Robson [14] for details).

Figure 2 (a) represents a dominance solvable game, for any payoff value  $\alpha \in \mathbf{R}$ . In the standard replicator dynamics (or any monotonic dynamics), the long-run outcome of this game has all players playing strategy 1. In contrast, if we introduce discriminating behavior by letting there be four types, where types 1 through 3 always play strategies 1 through 3, respectively,

while type 4 plays the best response to each of the four types (and therefore plays the Nash equilibrium strategy 1 against its own type), we can once again construct an augmented payoff-matrix which can be analyzed with the standard tools of evolutionary game theory (Banerjee and Weibull [1]).

$$\begin{pmatrix} 3 & -\alpha & 6 \\ 0 & -\alpha - 1 & 4 \\ 1 & -\alpha + 1 & 5 \end{pmatrix} \quad \begin{pmatrix} 3 & -\alpha & 6 & 3 \\ 0 & -\alpha - 1 & 4 & 4 \\ 1 & -\alpha + 1 & 5 & 1 \\ 3 & -\alpha + 1 & 6 & 3 \end{pmatrix}$$

Figure 2: A dominance-solvable game.

The augmented game in Figure 2 (b) has one connected set of Lyapunov (or weakly) stable states and one additional asymptotically stable state. The aggregate behavior in each population state in the first component is identical with the unique Nash equilibrium outcome of the original game (Fig.2 (a)), *i.e.*, all individuals use only strategy 1 and accordingly earn the Nash equilibrium payoff 3. Each such state is composed of individuals of types 1 and 4 only, and the share of the former exceeds 1/4.<sup>1</sup> In the second outcome, only types 2 and 4 survive, in shares 1/3 and 2/3, respectively, and all individuals earn the payoff  $(7 - \alpha)/3$ . Evidently, this outcome is *not* compatible with Nash equilibrium in the original game since the strictly dominated strategy 2 gets played. Moreover, the payoff in this asymptotically stable state depends on the value of  $\alpha$ , and for any  $\alpha$  exceeding  $-2/3$  this payoff falls below 3, the Nash equilibrium payoff. Thus, unlike in the previous example, what drives the population state away from the Nash outcome is not attraction towards some more efficient state.<sup>2</sup>

<sup>1</sup>To see that any state with no individuals of types 2 and 3, and with at least one quarter of type 1, is Lyapunov stable, note that on the boundary face of the mixed-strategy simplex where strategies 2 and 3 are extinct any state in the accordingly reduced game (with only strategies 1 and 4) is Lyapunov stable. Moreover in a neighborhood in the full mixed-strategy simplex of such a state, both strategies 2 and 3 earn less than 3, the payoff of 1 and 4, cf. Banerjee and Weibull (1991,1992).

<sup>2</sup>To see that the population state with no individuals of types 1 and 3, one-third of type 2 and two-thirds of type 4 is asymptotically stable, note that, when the dynamics is restricted to the corresponding boundary face of the mixed strategy simplex, this state is asymptotically stable. see e.g. Weibull (1992). Moreover, in a neighborhood of this

The last example we consider is the standard co-ordination game in Figure 3 (a). It has two asymptotically stable population states in the usual replicator dynamics, viz. the state in which all individuals play the “good” strategy 1 and the state in which all individuals play the “bad” strategy 2. However, if we identify types 1 and 2 with players who always play 1 and 2, respectively, and type 3 with players who play 1 against type 1, 2 against type 2 and 1 against type 3, then we obtain the augmented payoff matrix in Figure 3 (b). In this new game, strategy 2 (corresponding to individuals always playing 2 in the original game) cannot survive in the evolutionary selection process. Therefore, in the long run, in this setting with some discriminating behaviors, we necessarily end up in the “good” Nash equilibrium (1, 1) with payoff 2.

$$\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} \quad \begin{pmatrix} 2 & 0 & 2 \\ 0 & 1 & 1 \\ 2 & 1 & 2 \end{pmatrix}$$

Figure 3: A co-ordination game.

### 3 The model

The analysis in the present paper is restricted to the standard setting of evolutionary game theory, viz. finite and symmetric two-player games in normal form. Let  $I = \{1, 2, \dots, m\}$  be the set of *pure* strategies. Accordingly, a *mixed* strategy is a point  $x$  on the  $(m - 1)$ -dimensional unit simplex  $\Delta = \{x \in \mathbf{R}_+^m : \sum_i x_i = 1\}$  in  $m$ -dimensional Euclidean space. The *support* of a mixed strategy  $x \in \Delta$  is the subset  $C(x) = \{i \in I : x_i > 0\}$  of pure strategies which are assigned positive probabilities. A strategy  $x$  is called *interior* (or *completely mixed*) if  $C(x) = I$ , and we then write  $x \in \text{int}(\Delta)$ .

Let  $a_{ij}$  be the *payoff* of strategy  $i \in I$  when played against strategy  $j \in I$ , and let  $A$  be the associated  $k \times k$  payoff matrix. Accordingly, the (expected) payoff of a mixed strategy  $x \in \Delta$ , when played against a mixed strategy

---

state, both strategies 1 and 3 earn less than strategies 2 and 4, cf. Banerjee and Weibull (1991,1992).

$y \in \Delta$ , is  $u(x, y) = x \cdot Ay = \sum_{ij} x_i a_{ij} y_j$ . The payoff function  $u : \Delta^2 \rightarrow \mathbf{R}$  is bi-linear and the payoff of a pure strategy  $i \in I$ , when played against a mixed strategy  $y$ , is  $u(e^i, y)$ , where  $e^i$  is the  $i$ 'th unit vector in  $\mathbf{R}^m$ , etc. We will frequently identify a pure strategy  $i \in I$  with its mixed-strategy counterpart  $e^i \in \Delta$ . We summarize any symmetric 2-player normal-form game as a pair  $G = (I, u)$ , where  $I$  is the set of pure strategies and  $u : \Delta^2 \rightarrow \mathbf{R}$  the associated payoff function.

As usual, a pure strategy  $i \in I$  is *weakly dominated* if there exists a strategy  $x \in \Delta$  which never earns a lower payoff and sometimes a higher payoff (i.e.,  $u(x, y) \geq u(e^i, y) \quad \forall y \in \Delta$  with strict inequality for some  $y$ ). A pure strategy  $i \in I$  is *strictly dominated* if there exists a strategy  $x \in \Delta$  which always earns a higher payoff [ $u(x, y) > u(e^i, y) \quad \forall y \in \Delta$ ]. A *best reply* to a strategy  $y \in \Delta$  is a strategy  $x \in \Delta$  such that  $u(x, y) \geq u(x', y) \quad \forall x' \in \Delta$ . For each  $y \in \Delta$ , let  $\beta(y) \subset \Delta$  be its set of (mixed-strategy) best replies. A *Nash equilibrium* is a pair  $(x, y)$  of mutually best replies, a Nash equilibrium is *strict* if each strategy is the unique best reply to the other, and a Nash equilibrium  $(x, y)$  is *symmetric* if  $x = y$ . By Kakutani's Fixed Point Theorem, every finite and symmetric game has at least one symmetric Nash equilibrium.

One solution concept which is weaker than Nash equilibrium is iterative strict dominance. A pure strategy  $i \in I$  is said to be *iteratively strictly undominated* if it is not strictly dominated in the original game  $G$ , nor in the game  $G'$  obtained from  $G$  by removal of all strictly dominated strategies, nor in the game  $G''$  obtained from  $G'$  by removal of all strategies which are strictly dominated in  $G'$ , etc. A related but distinct solution concept is that of rationalizability (Bernheim [3]; Pearce [12]). A pure strategy  $i \in I$  is *never a best reply* if there exists no mixed strategy  $y \in \Delta$  against which  $i \in I$  is a best reply. A strategy  $i \in I$  is *rationalizable* if it is not a "never best reply" in the original game  $G$ , nor in the game  $G'$  obtained from  $G$  by removal of all "never best replies," nor in the game  $G''$  obtained from  $G'$  by removal of all "never best replies," etc. Each of these two methods of iterated elimination of pure strategies stops in a finite number of steps. Pearce [12] has shown that, while the two remaining sets may differ in games with more than two players, they in fact coincide in all two-player games.<sup>3</sup>

---

<sup>3</sup>A strictly dominated strategy is never a best reply, and hence the set of rationalizable strategies is always a subset of the set of strategies surviving the iterated elimination of strictly dominated strategies.

In this context of finite and symmetric two-person games in normal form, evolutionary game theory studies the long-run effects of evolutionary selection among strategies. Each individual is then (genetically or culturally) "programmed" to play some fixed strategy in the game in question, and individuals are randomly drawn in pairs from a large population to play the game. In the standard so-called *replicator dynamics*, the growth rate of the subpopulation programmed to a strategy  $i \in I$  is proportional to its current (average) payoff.

In this standard set-up of evolutionary game theory one can introduce a richer menu of behaviors as follows. Suppose there is a finite set  $T = \{1, \dots, k\}$  of *types* of individual. At each matching of two individuals from the population, both individuals know their own type and can (costlessly and without error) observe each other's type. One may then think of a possible *behavior* for an individual of any given type as a "rule"  $\varphi$  which to each type  $\tau \in T$  of opponent prescribes a strategy  $i = \varphi(\tau)$ , *i.e.* as a *function*  $\varphi : T \rightarrow I$ . Let  $F$  be the set of all such functions — the set of all possible behaviors — the number of which is  $\#F = m^k$ . We define a *character*  $\kappa$  as an element of the (finite) product set  $H = T \times F$ , *i.e.*, as a pair  $(\tau, \varphi)$ , where  $\tau \in T$  is a type, visible to each and everyone, and  $\varphi \in F$  is a behavior, observable only indirectly and incompletely via its outcome in interactions. The total number of possible characters thus is  $n = \#H = k \cdot m^k$ . Each individual in the population is fully described by her character since her own type determines the strategy choice of her "opponent" at each encounter, via the opponent's behavior, and her own behavior determines her strategy choice in every possible encounter. More precisely, the payoff to an individual with character  $\kappa = (\tau, \varphi)$ , when meeting an individual with character  $\lambda = (\nu, \psi)$ , is  $a_{ij}$  for  $i = \varphi(\nu)$  and  $j = \psi(\tau)$ . We will call a behavior  $\varphi \in F$  *constant* if there exists some  $i \in I$  such that  $\varphi(\nu) = i \forall \nu \in T$ .

This expanded setting for pairwise interactions is equivalent with the standard setting of evolutionary game theory if we take the pure strategy set to be  $H = \{1, 2, \dots, n\}$ , the set of characters, and the payoff matrix to be the  $n \times n$  matrix  $\mathcal{A}$  with entry  $\alpha_{\kappa\lambda} = a_{ij}$  in row  $\kappa$  and column  $\lambda$ , where  $\kappa = (\tau, \varphi)$ ,  $\lambda = (\nu, \psi)$ ,  $i = \varphi(\nu)$  and  $j = \psi(\tau)$ . In other words,  $\alpha_{\kappa\lambda}$  is the payoff that "strategy"  $\kappa$  earns when used against "strategy"  $\lambda$ . Since there are  $n$  pure strategies in this expanded game, its mixed-strategy space is the  $(n-1)$ -dimensional unit simplex  $\Sigma$  in  $\mathbf{R}^n$ . For any pair of "mixed strategies"  $\sigma, \mu \in \Sigma$  in the expanded game, let  $v(\sigma, \mu)$  denote the payoff to "strategy"

$\sigma$  when used against “strategy”  $\mu$ , *i.e.*,

$$v(\sigma, \mu) = \sigma \cdot \mathcal{A} \mu = \sum_{\kappa, \lambda \in H} \sigma_{\kappa} \cdot \alpha_{\kappa\lambda} \cdot \mu_{\lambda},$$

In this way, we have defined the payoff function  $v : \Sigma \times \Sigma \rightarrow \mathbf{R}$  of the expanded game  $\mathcal{G} = (H, v)$  in which characters are viewed as pure strategies. (The standard setting of evolutionary game theory can be identified with the special case when  $T$  is a singleton set, *i.e.*, when there is only one type  $\tau$ . For then there are precisely  $k$  behaviors, *viz.*  $\varphi(\tau) = i$ , for  $i = 1, 2, \dots, k$ . Hence,  $H = I$ ,  $\Sigma = \Delta$  and  $v = u$ .)

We need some notation to describe the composition and dynamics of the population in the expanded game  $\mathcal{G}$ . For  $\kappa \in H$ , let  $p_{\kappa}$  be the population share of individuals with character  $\kappa$ . The vector  $p = (p_{\kappa})_{\kappa \in H}$  is then the population *state*, a point on the unit simplex  $\Sigma$ , the space of mixed strategies in  $\mathcal{G}$ , and one may study the workings of the (standard, continuous-time) replicator dynamics on this (typically very high-dimensional) simplex. When the population state is  $p \in \Sigma$ , any individual of character  $\kappa \in H$ , *i.e.*, playing “strategy”  $\kappa$  in game  $\mathcal{G}$ , obtains the (average, expected) payoff  $v(e^{\kappa}, p)$  at a random matching, where  $e^{\kappa}$  is the  $\kappa$ :*th* unit vector in  $\mathbf{R}^n$ . Likewise, the average payoff in the population is  $v(p, p)$ . The payoff function  $v$  is bi-linear and the replicator dynamics is

$$\dot{p}_{\kappa} = v(e^{\kappa} - p, p)p_{\kappa} \quad [\forall \kappa \in H],$$

to which all results from standard evolutionary game theory apply.

Applying this general framework to the Prisoner’s Dilemma game in Figure 1 (a), we note that the partially expanded game in Figure 1 (b) has three types,  $T = \{1, 2, 3\}$ , but each type is restricted to a single behavior. For in that example, the only behavior of type 1 is the constant behavior  $\varphi(\nu) = 1$  for all  $\nu \in T$ , the only behavior of type 2 is likewise constant,  $\psi(\nu) = 2$  for all  $\nu \in T$ , and the only behavior of type 3 is the non-constant behavior  $\gamma(1) = \gamma(2) = 2$  and  $\gamma(3) = 1$ . In contrast, the fully expanded game  $\mathcal{G}$  has  $2^3 = 8$  behaviors available to each of the three types, so the game  $\mathcal{G}$  has  $8 \times 3 = 24$  pure strategies (characters). It is easily verified that, unlike in the partially expanded game in Figure 1 (b), the strategy  $\kappa = (\tau, \psi)$ , which (for any given type  $\tau$ ) plays 2 against every type  $\nu \in T$ , is *not* weakly dominated in the fully expanded game  $\mathcal{G}$ . Similarly, if we construct the fully

expanded game  $\mathcal{G}$  associated with the co-ordination game in Figure 3 (a), then play of the “bad” Nash equilibrium strategy 2 against all types is not weakly dominated (unlike in the partially expanded game in Figure 3 (b)).

## 4 Basic results

For any type  $\tau \in T$ , we will call the subset of individuals of type  $\tau$  *sub-population*  $\tau$ . Hence, individuals in the same sub-population differ only with respect to their behaviors. In particular, they are all met by the same choice of strategy when meeting any given individual; if an individual of type  $\tau$  meets an individual with behavior  $\varphi$ , then the latter will choose strategy  $\varphi(\tau) \in I$ . Let  $p^\tau$  denote the population share of individuals of type  $\tau$ , *i.e.*, for any population state  $p \in \Sigma$  and type  $\tau \in T$ ,  $p^\tau$  is the sum of all  $p_\kappa$  such that  $\kappa = (\tau, \varphi)$  for some  $\varphi \in F$ .

It turns out to be useful to decompose aggregate behavior in the whole population into the aggregate behaviors in the matchings between every combination of sub-populations  $\tau, \nu \in T$  (including the case  $\tau = \nu$ ). For each pure strategy  $i \in I$  in the original game  $G$ , every type  $\tau \in T$  and population state  $p \in \Sigma$  with  $p^\tau > 0$ , let  $p_i^{\tau\nu}$  be the share of individuals in sub-population  $\tau$  who use strategy  $i$  when meeting an individual of type  $\nu \in T$ , *i.e.*,  $p_i^{\tau\nu}$  is the share of individuals of type  $\tau$  who have behaviors  $\varphi$  with  $\varphi(\nu) = i$ . Clearly the vector  $p^{\tau\nu} = (p_i^{\tau\nu})_{i \in I}$  is a point on the unit simplex  $\Delta$  of the original game  $G$ , *i.e.*  $p^{\tau\nu}$  can be viewed as the mixed strategy (in  $G$ ) that any individual of type  $\nu$  faces when matched with an individual of type  $\tau$ . In other words,  $p^{\tau\nu} \in \Delta$  is the aggregate behavior of individuals in sub-population  $\tau \in T$  when meeting individuals from sub-population  $\nu \in T$ , and  $p^{\nu\tau} \in \Delta$  is the aggregate behavior of sub-population  $\nu$  in the same encounters.

In the replicator dynamics, as applied to the standard setting of (single-population) evolutionary game theory, any Lyapunov stable population state  $x \in \Delta$ , viewed as a mixed strategy, is a best reply to itself, *i.e.*,  $(x, x)$  constitutes a Nash equilibrium (Bomze [6]).<sup>4</sup> This fundamental result generalizes to the present setting as follows: if sub-populations  $\tau$  and  $\nu$  (possibly

---

<sup>4</sup>Note that Lyapunov (or weak) stability is less stringent than the perhaps more familiar criterion of asymptotic stability: Lyapunov stability essentially prohibits local “drift” away from the stationary state in question, while asymptotic stability essentially requires a local “pull” towards the state, see e.g. Hirsch and Smale (1974) for definitions.

$\tau = \nu$ ) are non-extinct in a Lyapunov stable population state, then these sub-populations play some Nash equilibrium against each other. Formally:

**Proposition 1** *If  $p \in \Sigma$  is Lyapunov stable in the replicator dynamics on  $\Sigma$  and  $p^\tau, p^\nu > 0$ , then  $(p^{\tau\nu}, p^{\nu\tau}) \in \Delta^2$  is a Nash equilibrium of  $G$ .*

**Proof:** Suppose  $p \in \Sigma$  is stationary with  $p^\tau, p^\nu > 0$ , and suppose  $p^{\tau\nu} \in \Delta$  is not a best reply (in  $G$ ) to  $p^{\nu\tau} \in \Delta$ . Then some strategy  $i \in I$  in the support  $C(p^{\tau\nu}) \subset I$  earns a suboptimal payoff against  $p^{\nu\tau}$ . Let  $\lambda \in H$  be any character  $(\tau, \varphi)$  with  $\varphi(\nu) = i$  and  $p_\lambda > 0$  (such a  $\lambda$  exists since  $p^\tau > 0$  and  $p_i^{\tau\nu} > 0$ ). Let  $\kappa \in H$  be a character  $(\tau, \psi)$ , where  $\psi(\nu) \in \beta(q^{\nu\tau})$  and  $\psi(\omega) = \varphi(\omega) \quad \forall \omega \in T, \omega \neq \nu$ . In other words, character  $\kappa$  plays a best reply against the type- $\nu$  sub-population, hence earning a higher payoff than character  $\lambda$  in such encounters, and otherwise  $\kappa$  plays exactly like  $\lambda$ . Since  $p^\nu > 0$ , we thus have  $v(e^\kappa - e^\lambda, p) > 0$ . By stationarity of  $p_\lambda > 0$ ,  $v(e^\lambda - p, p) = 0$ , so  $v(e^\kappa - p, p) > 0$ , implying  $p_\kappa = 0$  by stationarity. However, by continuity of  $v$ ,  $\dot{q}_\kappa = v(e^\kappa - q, q)q_\kappa > 0$  for all interior population states  $q$  in some neighborhood of  $p$ . Hence,  $p$  is not Lyapunov stable in the replicator dynamics on  $\Sigma$ . ■

In other words, even when allowing for all possible behaviors, and not just fixed strategies, evolution selects, in so far as (Lyapunov) stable outcomes are concerned, behavior which is "rational" and "coordinated" in the sense of Nash-equilibrium play between all sub-populations. As a consequence, (Lyapunov) stable aggregate behavior is always some convex combination of Nash equilibrium play. That the converse is not generally true, *i.e.*, that certain convex combinations of Nash equilibria may be dynamically unstable, is seen in the co-ordination game in Figure 3 (a) above. The mixed-strategy Nash equilibrium strategy  $x \in \Delta$  in this game is to randomize with probabilities 1/3 and 2/3, respectively, and it is well-known from the literature on standard evolutionary games that this state is unstable in the replicator dynamics. (To see this, let  $\sigma$  be the stationary state in which the population share playing strategy 1 is 1/3 and the population share playing strategy 2 is 2/3. Then a slightly higher share of players using strategy 1 gives these a higher than average payoff, and so their population share increases, implying the instability of  $\sigma$ .) As a consequence, this state is unstable also in the present setting, since with only one type,  $T = \{1\}$ , we are back in the traditional setting.

A related result for standard evolutionary games is that if an interior dynamic solution path  $x(t)$  to the replicator dynamics converges over time, then the limit state  $y \in \Delta$  is a best reply to itself, *i.e.*, then  $(y, y)$  constitutes a Nash equilibrium of  $G$  (Nachbar [11]). The generalization of this result is straight-forward: if an interior dynamic solution path  $p(t)$  to the replicator dynamics on  $\Sigma$  converges over time to some state  $q \in \Sigma$ , and if sub-populations  $\tau$  and  $\nu$  (possibly  $\tau = \nu$ ) are non-extinct in that limiting state, then they play a Nash equilibrium against each other. Formally:

**Proposition 2** *If an interior dynamic solution path to the replicator dynamics converges to a state  $p \in \Sigma$  with  $p^\tau, p^\nu > 0$ , then  $(p^{\tau\nu}, p^{\nu\tau}) \in \Delta^2$  is a Nash equilibrium of  $G$ .*

**Proof:** Suppose  $p(0) \in \text{int}(\Sigma)$ ,  $p(t) \rightarrow p$ , and  $p^\tau, p^\nu > 0$ . Suppose  $p^{\tau\nu} \in \Delta$  is *not* a best reply (in  $G$ ) to  $p^{\nu\tau} \in \Delta$ . Then some strategy  $i \in I$  in the support  $C(p^{\tau\nu}) \subset I$  earns a suboptimal payoff against  $p^{\nu\tau}$ . Just as in the proof of Prop.1, let  $\lambda \in H$  be any character  $(\tau, \varphi)$  with  $\varphi(\nu) = i$  and  $p_\lambda > 0$ , and let  $\kappa \in H$  be a character  $(\tau, \psi)$ , where  $\psi(\nu) \in \beta(p^{\nu\tau})$  and  $\psi(\omega) = \varphi(\omega) \forall \omega \in T$  with  $\omega \neq \nu$ . Since  $p^\nu > 0$ , we thus have  $v(e^\kappa - e^\lambda, p) > 0$ . The limiting state  $p$  is stationary by continuity of  $v$ , and  $p^\lambda > 0$ , so  $v(e^\lambda - p, p) = 0$ . Hence,  $v(e^\kappa - p, p) > 0$ , implying  $p^\kappa = 0$  by stationarity. By continuity of  $v$ ,  $\dot{q}_\kappa = v(e^\kappa - q, q)q_\kappa > 0$  for all states  $q \in \text{int}(\Sigma)$  in some neighborhood of  $p$ . Thus the coordinate  $p_\kappa(t)$  does not converge to zero along the path  $p(t)$ , a contradiction to  $p_\kappa = 0$ . Hence,  $p^{\tau\nu} \in \beta(p^{\nu\tau})$ . ■

Samuelson and Zhang [16] have shown that the replicator dynamics, in the standard (two-population) setting of evolutionary game theory, wipe out all non-rationalizable strategies, irrespective of whether the population state converges or not. More exactly, if a strategy  $i \in I$  is not rationalizable, then its population share converges to zero in the continuous-time replicator dynamics on  $\Delta$ , from any interior initial state. In particular, all strictly dominated strategies vanish. However, the same is not true for weakly dominated strategies, an issue addressed in Samuelson [15], where it is shown that a variety of evolutionary selection dynamics, including the replicator dynamics, do not eliminate such strategies. This observation has an important implication for the present generalized setting, since if there is more than one type in the population, then any strategy  $i \in I$  which is strictly dominated in the orig-

inal game  $G$  is only weakly dominated in the expanded game  $\mathcal{G}$ .<sup>5</sup> Therefore we cannot invoke Samuelson's and Zhang's [16] result to prove that strategies  $i \in I$  which are strictly dominated in  $G$  will be wiped out in the replicator dynamics on the strategy space  $\Sigma$  of the expanded game  $\mathcal{G}$ . However, in all examples in [15], the survival of weakly dominated strategies is due to the fact that the strategies against which they fare badly vanish. Indeed, one can show that if a pure strategy  $i \in I$  is weakly dominated by some strategy  $m \in \Delta$ , and sub-population  $i$  does not vanish over time, then it must be the case that all sub-populations  $j$  against which  $m$  is better than  $i$  vanish (see Appendix) In force of this result, we have

**Proposition 3** *If strategy  $i \in I$  is non-rationalizable in  $G$ , then the product  $p_i^{\tau\nu} \cdot p^\nu$  converges to zero over time along any interior solution path to the replicator dynamics,  $\forall \tau, \nu \in T$ .*

**Proof:** Suppose first that strategy  $i \in I$  is strictly dominated by  $m \in \Delta$  in  $G$ . For any types  $\tau, \nu \in T$ , let  $H_i^{\tau\nu} = \{\kappa \in H : \kappa = (\tau, \varphi) \text{ for some } \varphi \in F \text{ such that } \varphi(\nu) = i\}$  and  $H^\nu = \{\lambda \in H : \lambda = (\nu, \psi) \text{ for some } \psi \in F\}$ . Then  $p_i^{\tau\nu}$  is the sum of all  $p_\kappa$  with  $\kappa \in H_i^{\tau\nu}$  and  $p^\nu$  is the sum of all  $p_\lambda$  with  $\lambda \in H^\nu$ . From now on, fix any character  $\kappa = (\tau, \varphi) \in H_i^{\tau\nu}$  and let  $\sigma_\kappa \in \Sigma$  be such that, for each  $j \in C(m)$ , the sum of  $\sigma_\kappa(\tau, \psi)$ , over all  $\psi \in F$  such that  $\psi(\nu) = j$  and  $\psi(\rho) = \varphi(\rho) \forall \rho \neq \nu$ , is  $m_j$ . In other words, the mixed strategy  $\sigma_\kappa$  in the expanded game  $\mathcal{G}$  assigns zero probability to all pure strategies  $\lambda = (\rho, \psi) \in H$  which have types  $\rho \neq \tau$  and/or have behaviors  $\psi$  which differ from  $\varphi$  when meeting individuals of types  $\rho \neq \nu$  and/or have behaviors  $\psi$  which play strategies  $j \in I$  outside the support of  $m \in \Delta$ . Moreover,  $\sigma_\kappa$  randomizes the behavior against type  $\nu$  in such a way that the induced distribution over the pure strategy set  $I$  coincides with that of  $m \in \Delta$ . It follows that  $\kappa \in H$  is weakly dominated by  $\sigma_\kappa \in \Sigma$  (because  $\kappa$  does worse than  $\sigma_\kappa$  against pure strategies  $\lambda \in H^\nu$  and does equally well against all other pure strategies  $\lambda \in H$ ). Formally,  $v(\sigma_\kappa - e^\kappa, e^\lambda) > 0 \forall \lambda \in H^\nu$  and  $v(\sigma_\kappa - e^\kappa, e^\lambda) = 0 \forall \lambda \notin H^\nu$ . By the proposition in the Appendix,  $p_\kappa(t)p_\lambda(t) \rightarrow 0$  as  $t \rightarrow \infty, \forall \lambda \in H^\nu$ . Summing over all  $\lambda \in H^\nu$ , we get  $p_\kappa(t)p^\nu(t) \rightarrow 0$ . This

<sup>5</sup>This is so because in the present setting a pure strategy,  $i$ , in the underlying game is only part of the pure strategies in the expanded game. In particular, there are individuals who in the expanded game use strategy  $i$  against some opponents and other strategies against others.

is true for any  $\kappa \in H_i^{\tau\nu}$ , so we may also sum over all  $\kappa \in H_i^{\tau\nu}$ , yielding  $p_i^{\tau\nu}(t)p^\nu(t) \rightarrow 0$  as  $t \rightarrow +\infty$ .

It follows that after a sufficiently long time interval, the game will be played arbitrarily close to how it would have been played, had every strictly dominated strategy been eliminated from game  $G$ . Hence, if a strategy  $i \in I$  is not strictly dominated in  $G$  but in the game  $G'$  obtained when all strictly dominated strategies have been eliminated, then the above logic implies, by continuity of  $v$ , that the product  $p_i^{\tau\nu}(t)p^\nu(t)$  vanishes as  $t \rightarrow +\infty$ . Indeed, this result holds for all strategies  $i \in I$  which do not survive the iterated elimination of strictly dominated strategies, a criterion which is equivalent to being non-rationalizable in two-person games, thereby establishing the claim in the proposition. ■

The above result says, inter alia, that the examples mentioned above (Robson [14]; Banerjee and Weibull [1] and Banerjee and Weibull [2]) in an essential way depend on the incompleteness of the population. To see this, let  $z_i$  be the share of *matchings* at which strategy  $i \in I$  is used in any interior population state  $p \in \Sigma$  :

$$z_i = \sum_{\tau \in T} p^\tau \left[ \sum_{\nu \in T} p_i^{\tau\nu} p^\nu \right].$$

Hence  $z = (z_1, \dots, z_k) \in \Delta$  and any interior dynamic path  $p(t)$  on  $\Sigma$  induces an associated dynamic path  $z(t)$  on  $\Delta$ . The following result is an immediate implication of Proposition 3: if strategy  $i \in I$  is non-rationalizable in  $G$ , then  $z_i(t) \rightarrow 0$  as  $t \rightarrow \infty$ , along any interior solution path  $p(t)$  to the replicator dynamics.

## 5 Limit results

We here investigate the limiting properties of the set of payoffs in stationary and Lyapunov stable states, respectively, as the number of types tends to infinity. After all, it is usually not a priori evident how many types one should naturally assume in a given modelling context, and even the set  $T$  of types could itself follow some evolutionary process over time. Our first result establishes that the set of stationary payoffs expands towards a limit set  $V_\infty$  which is a dense subset of the set  $U^{NE}$  of "symmetric convex combinations" of Nash equilibrium payoffs in the underlying game  $G$ . The second

result establishes that the subset of Lyapunov stable stationary payoffs has a limit set  $W_\infty \subset V_\infty$  which is a proper, not dense, subset of the set  $U^{NE}$ .

In order to make precise these claims, some more notation is needed. First, let  $NE(G) \subset \Delta^2$  be the set of Nash equilibria of  $G$ . Recall that there exists at least one symmetric Nash equilibrium. Let  $(m, m) \in \Delta^2$  be such. Let  $W^{NE} \subset \mathbf{R}^2$  be the set of Nash equilibrium payoff pairs, and let  $D \subset \mathbf{R}^2$  be the diagonal in  $\mathbf{R}^2$  (i.e.,  $D = \{(x, x) : x \in \mathbf{R}\}$ ). The set  $W^{NE}$  is compact, symmetric around  $D$ , and, as just noted, contains at least one point  $(u^m, u^m)$  in  $D$  (where  $u^m = u(m, m)$ ). Let  $W^{SNE}$  be the (nonempty) intersection of the convex hull of  $W^{NE}$  with  $D$ , i.e.,  $W^{SNE} = co(W^{NE}) \cap D$ . Then  $W^{SNE} = D \cap [\underline{u}, \bar{u}]$  for some  $\underline{u}, \bar{u} \in \mathbf{R}$  with  $\underline{u} \leq u^m \leq \bar{u}$ . Let  $U^{NE} = [\underline{u}, \bar{u}]$ , i.e.,  $U^{NE}$  is the set of "symmetric convex combinations" of Nash equilibrium payoffs mentioned above.

Turning to the expanded game  $\mathcal{G}$ , we first observe that stationarity in the replicator dynamics requires all individuals in the population to earn the same payoff. Let  $V \subset \mathbf{R}$  be the set of payoffs compatible with stationarity in a given expanded game  $\mathcal{G}$ . (More exactly,  $w \in V$  iff  $w = v(p, p)$  for some stationary state  $p \in \Sigma$ .) Clearly the set  $V$  depends on the set  $T$  of types only via the cardinality of  $T$ ; if  $\#S = \#T$  then the associated sets  $V$  coincide. Hence, without ambiguity we may write  $V_k \subset \mathbf{R}$  for the set of stationary payoffs in any expanded game  $\mathcal{G}$  of the same underlying game  $G$ , with  $k$  types.

**Proposition 4** *The sequence of payoff sets  $V_k$  is (weakly) increasing towards a limit set  $V_\infty$  which is dense in  $U^{NE}$ .*

**Proof:** It is evident that the sequence of sets  $V_k$  is (weakly) increasing, since, when going from  $k$  to  $k+1$  types one can always let the  $k$  "old" types be present in the same population shares and "behave" as before and let the population share of the "new" type be zero. Moreover, all sets  $V_k$  are subsets of the image of the compact set  $\Delta^2$  under the continuous function  $u$ , and so the sequence has a limit ("its smallest upper bound"), which we denote  $V_\infty$ . Clearly  $u^m \in U^{NE}$  belongs to each set  $V_k$ ; just let every type use only constant behaviors such that the population shares across strategies is identical with the Nash equilibrium strategy  $m \in \Delta$ . Then all individuals use pure strategies in the support of  $m$  and meet  $m$ , and so earn the same (average) payoff, implying stationarity. Now consider any point  $w$  between

$u^m$  and  $\bar{u}$ , or between  $\underline{u}$  and  $u^m$ . Then  $w$  is a convex combination of  $u^m$  and  $\bar{u}$  or  $\underline{u}$ , with weight  $\lambda \in (0, 1)$  attached to  $\bar{u}$  or  $\underline{u}$ , as the case may be. By having sufficiently many types in the expanded game, one can obtain a stationary payoff corresponding to a rational  $\hat{\lambda} \in (0, 1)$  arbitrarily close to  $\lambda$ , thus establishing the denseness of  $V_\infty$  in  $U^{NE}$ .

The logic is the same in both cases, so suppose for the sake of definiteness that it is  $\bar{u}$  which is the second payoff. By definition of  $U^{NE}$ , there exist strategies  $x$  and  $y$  in  $\Delta$  such that  $(x, y), (y, x) \in NE(G)$  and  $\bar{u} = [u(x, y) + u(y, x)]/2$ . For any  $\epsilon > 0$  there exists some  $\hat{\lambda} = s/k \in [0, 1]$  within distance  $\epsilon$  from  $\lambda$  such that  $s$  is even and  $k \geq s + 1$ . Let there be  $k$  types in  $T$ , and place all types around a circle. Let each type  $\tau \in T$  play  $m \in \Delta$  against its own type,  $x \in \Delta$  against its  $s/2$  nearest "clockwise neighbor" types on the circle,  $y \in \Delta$  against its  $s/2$  nearest "counter-clockwise" neighbors, and  $m$  against all  $k - s - 1$  other types on the circle. Then all types play Nash equilibria with each other, and all individuals in the population earn the same (average) payoff  $[1 \cdot u^m + \frac{s}{2} \cdot u(x, y) + \frac{s}{2} \cdot u(y, x) + (k - s - 1) \cdot u^m]/k = \hat{\lambda} \cdot u_m + (1 - \hat{\lambda}) \cdot \bar{u}$ . ■

The finding (Proposition 1) that Lyapunov stability implies Nash equilibrium play within each sub-population  $\tau \in T$ , as well as between any two sub-populations  $\tau, \nu \in T$ , has the immediate implication that all individuals in any stationary Lyapunov stable state earn the same payoff and that this payoff belongs to  $U^{NE}$ . Let  $W_k \subset U^{NE}$  denote the subset of payoffs compatible with Lyapunov stability in an expanded game  $\mathcal{G}$  with  $k$  types. (i.e.,  $w \in W_k$  iff  $w = v(p, p)$  for *some* Lyapunov stable state  $p \in \Sigma$ .) The following result establishes that, as the number  $k$  of types increases towards infinity, the Lyapunov stable payoff set  $W_k$  converges towards a limit set  $W^\infty \subset U^{NE}$  (in a sense to be made precise below). Moreover, the limit set  $W^\infty$  may be a *proper* subset of  $U^{NE}$  in the strong sense that it differs from  $U^{NE}$  by at least some interval of positive length. (In particular,  $W^\infty$  is not necessarily dense in  $U^{NE}$ .)

In order to establish this claim, we treat the possibility of minimax Nash equilibria separately. For this purpose, let  $u_o \in \mathbf{R}$  be the *minimax payoff* in  $G$  and let  $\Delta_o \subset \Delta$  be the (nonempty) set of *minimax strategies* in  $G$ , i.e.,

$$u_o = \min_{y \in \Delta} \max_{x \in \Delta} u(x, y)$$

and  $\Delta_o = \{y \in \Delta : u(x, y) \leq u_o \ \forall x \in \Delta\}$ . Clearly no Nash equilibrium payoff is lower than the minimax value:  $u_o \leq \underline{u}$ . Let  $W_k^+ = W_k \cap (u_o, +\infty)$ .

**Proposition 5** *The sequence of payoff sets  $W_k^+$  is (weakly) increasing towards a limit set in  $U^{NE}$ . If there is at most one Nash equilibrium in  $G$  with payoffs  $(u_o, u_o)$ , then  $W_k \rightarrow W^\infty \subset U^{NE}$ . There are games  $G$  such that  $U^{NE}$  contains an interval disjoint from  $W^\infty$ .*

**Proof:** We first establish  $W_k^+ \subset W_{k+1}^+ \forall k$ , where each set  $W_k^+ \subset U^{NE}$  is a subset of the image of the compact set  $\Delta^2$  under the continuous mapping  $u$ , and hence the sequence  $(W_k^+)_{k \rightarrow \infty}$  has a compact upper bound and thus converges to some subset of  $U^{NE}$  if the sequence  $\{W_k^+\}$  is (weakly) increasing.

For any number  $k$  of types, let  $H_k$  be the set of characters,  $F_k$  the set of behaviors, and  $\Sigma_k$  the unit simplex, in the associated expanded game  $\mathcal{G}$ . Suppose  $p \in \Sigma_k$  is Lyapunov stable in the replicator dynamics on  $\Sigma_k$  and  $w = v(p, p) > u_o$ . Then every  $\kappa \in H_k$  with  $p_\kappa > 0$  earns  $w$ , i.e.,  $v(e^\kappa, p) = w$ , by stationarity. Given such a state  $p \in \Sigma_k$  one may identify an associated state  $q \in \Sigma_{k+1}$  in the game  $\mathcal{G}'$  with one more type  $\tau'$ , a state which is Lyapunov stable in the replicator dynamics on  $\Sigma_{k+1}$  and has average payoff  $w$ , as follows. When this is done, the claim  $W_k^+ \subset W_{k+1}^+$  has been established.

For this purpose, let  $m \in \Delta_o$  be a minimax strategy in  $G$ , and let  $T' = T \cup \{\tau'\}$ , where  $\#T = k$  and  $\tau' \notin T$ . We say that a character  $\lambda = (\nu, \psi) \in H_{k+1}$  agrees with a character  $\kappa = (\tau, \varphi) \in H_k$  if  $\nu = \tau$  and  $\psi(\nu) = \varphi(\nu) \forall \nu \in T$ . For each  $\lambda = (\nu, \psi) \in H_{k+1}$  with  $\nu \in T$  there exists precisely one  $\kappa = (\tau, \varphi) \in H_k$  with which  $\lambda$  agrees. Denote this  $\kappa(\lambda)$ . Define the state  $q \in \Sigma_{k+1}$  by letting, for each  $\lambda = (\nu, \psi) \in H_{k+1}$  with  $\nu \in T$ ,  $q_\lambda = m_i p_{\kappa(\lambda)}$  if  $\psi(\tau') = i$ , and, for each  $\lambda = (\nu, \psi) \in H_{k+1}$  with  $\nu \notin T$ ,  $q_\lambda = 0$ . In other words, sub-population  $\tau'$  is extinct in state  $q \in \Sigma_{k+1}$  and all other sub-populations behave against each other precisely as in state  $p \in \Sigma_k$  and they all minimax individuals of the "new" type  $\tau'$ , if these would appear. Hence, the average payoff  $v(q, q)$  in state  $q \in \Sigma_{k+1}$  is  $w$ . Moreover,  $q$  is Lyapunov stable. To see this, note that any "pure strategy"  $\lambda \in H_{k+1}$  of type  $\tau'$  earns less than  $w$  in state  $q$ :  $v(e^\lambda - q, q) \leq u_o - w < 0 \forall \lambda \in H^{\tau'}$ . By continuity of  $v$ , this inequality holds for all  $p \in \Sigma_{k+1}$  in some neighborhood of  $q$ . Thus,  $\dot{p}_\lambda < 0$  for all  $\lambda \in H^{\tau'}$  in this neighborhood. Since by hypothesis the state  $q \in \Sigma_{k+1}$  is Lyapunov stable relative to the boundary face where  $q_\lambda = 0, \forall \lambda \in H^{\tau'}$ ,  $q$  is Lyapunov stable in the full replicator dynamics on  $\Sigma_{k+1}$ .

If  $u_o \notin U^{NE}$ , then, by Prop.1,  $W_k^+ = W_k$  for all  $k$ , since then  $(u_o, u_o) \notin W^{NE}$ , and so the limit  $W^\infty$  exists. If, on the other hand,  $u_o \in U^{NE}$ , then  $u_o = u(m, m)$  for some Nash equilibrium  $(m, m) \in \Delta_o$ , and  $(u_o, u_o)$  cannot

be obtained as the convex combination of Nash equilibria with other payoffs. Suppose  $u_o \notin W_k$  for some  $k$ , i.e., every corresponding population state  $p \in \Sigma_k$  is unstable, and suppose also  $u_o \in W_{k+1}$ , i.e., some corresponding population state  $q \in \Sigma_{k+1}$  is Lyapunov stable. However, by stationarity, all individuals in  $q$  earn  $u_o$ , and, by Prop.1, all present types play some minimax Nash equilibrium strategy against each present type. If the number of present types is less than  $k + 1$  we have a contradiction to the hypothesis that all corresponding population states  $p \in \Sigma_k$  are unstable. If the number of present types is  $k + 1$ , we again have a contradiction, since if every type plays some minimax Nash equilibrium strategy against each type, we could take away one type and re-scale all  $k$  remaining subpopulation distributions accordingly, without losing stability, thus creating a Lyapunov stable state  $p \in \Sigma_k$  with payoff  $u_o$ .<sup>6</sup>

The possibility of the existence of an interval  $(a, b) \subset \mathbf{R}$  which is disjoint from  $W^\infty$  and yet is a subset of  $U^{NE}$  follows from the example given below. ■

Consider the co-ordination game in Figure 3 (a) played by two types,  $T = \{1, 2\}$ . Let  $p \in \Sigma$  be a stationary state with a payoff slightly above 1. Such a state  $p$  can be generated in only two distinct ways: either at least one type plays the "good" Nash-equilibrium strategy 1 against itself, or each type plays this strategy when meeting the other type. In the first case, the type playing the "good" equilibrium strategy with itself can "invade" the population, and thus destabilize the state  $p$ . In the second case, a small invasion by a population consisting of 50% of each type can "invade." Therefore, no stationary state which has a payoff slightly above 1 is stable. In contrast, both the "good" (strict) Nash equilibrium payoff 2 and the "bad" (strict) Nash equilibrium payoff 1 are always Lyapunov stable, irrespective of the number of types. This follows immediately from the well-known fact that every strict equilibrium is an asymptotically stable population state in standard evolutionary dynamics (see e.g. Samuelson and Zhang [16]), combined with the monotonicity property established in Proposition 5.

However, there are games in which a unique Pareto dominant (but non-strict) Nash equilibrium is unstable. An example is given in Figure 4. Here  $\alpha \in (0, 1)$ ,  $\beta \in (0, 1 - \frac{\alpha}{3})$  and  $\gamma \geq 0$ . The three first rows and columns to-

---

<sup>6</sup>The so-called Hawk-Dove game is an example of a game in which a minimax Nash equilibrium payoff is stable in the case of one type and unstable in the case of two or more types. A demonstration of this is available from the authors upon request.

gether constitute a generalized Rock-Paper-Scissors game which has a unique Nash equilibrium, and in this equilibrium both players randomize uniformly over the three strategies, and each player obtains the payoff  $1 - \frac{\alpha}{3}$ . It is well-known that this equilibrium is unstable in the usual single-population replicator dynamics (see e.g. Hofbauer and Sigmund [9] or Weibull [20]). For nonnegative values of  $\gamma$ , this "Rock-Papers-Scissors equilibrium" remains a Nash equilibrium in the full game in Figure 4. However, the full game has two more Nash equilibria, each of which is symmetric. One is the strict equilibrium in which both players use only strategy 4, resulting in payoff  $\beta$  to both players - by hypothesis a lower payoff than in the "Rock-Papers-Scissors equilibrium". The third Nash equilibrium is completely mixed and its payoff can be made arbitrarily low by choosing  $\gamma$  sufficiently large. However, the unique Pareto-dominant Nash equilibrium, giving payoff  $1 - \frac{\alpha}{3}$  to each player, is not Lyapunov stable, for any number of types. For if  $p \in \Sigma$  is a stationary population state with this payoff, then all individuals in the population earn the same payoff and all types play the "Rock-Papers-Scissors equilibrium" with all types. When a type meets itself, the situation is exactly as in the usual single-population case, and thus the state is unstable with respect to pairwise matchings within each type. When two different types meet, the situation is the same as in the two-population case (cf. Hofbauer and Sigmund [9]). But instability in the single-population dynamics implies instability in the standard two-population dynamics. Thus, no stationary state which produces the Pareto-dominant Nash equilibrium payoff is Lyapunov stable in this example.

$$\begin{pmatrix} 1 & 2 - \alpha & 0 & -\gamma \\ 0 & 1 & 2 - \alpha & -\gamma \\ 2 - \alpha & 0 & 1 & -\gamma \\ -\gamma & -\gamma & -\gamma & \beta \end{pmatrix}$$

**Figure 4:** An enlarged "Rock-Paper-Scissors" game.

## 6 Concluding discussion

The results presented in the preceding sections evidently depend, *inter alia*, on our choice of solution concept. For the most part, we have used Lyapunov (or weak) stability as a criterion for selection among stationary population states in our evolutionary model of discriminating behaviors. This dynamic criterion is weaker than the maybe more familiar criterion of asymptotic stability which in essence requires that the state have a basin of attraction which contains some neighborhood of the state. For unlike asymptotic stability, Lyapunov stability allows for the absence of a local “pull” back towards the state in question, so such a state may even have an empty basin of attraction.

Hence, one might wish to use a stricter dynamic solution criterion than Lyapunov stability. However, asymptotic stability has too much cutting power in the present context, since there may be many stationary states in which several types of player occupy exactly the same role in the sense of having the same behavior (rule to select a strategy against an opponent) and receiving the same responses from opponents. As a result, the population shares of these types can be altered without altering the outcome, implying that such stationary states cannot be asymptotically stable.

An interesting alternative dynamic solution concept is *set-wise asymptotic stability*. In essence, a set of stationary states is asymptotically stable if it has a basin of attraction which contains a neighborhood of the set. (See Swinkels [18] and Ritzberg and Weibull [13] for applications of such criteria to standard evolutionary game theory.) Applying this criterion to the co-ordination game of Section 2, it turns out that the (Lyapunov stable) “bad” (strict) Nash outcome can be eliminated. To see why this is the case, suppose  $P \subset \Sigma$  is an asymptotically stable set of stationary states yielding the inefficient Nash equilibrium payoff 1, then all individuals have to play the “bad” strategy 2 against each other. If the set  $P$  is nonempty, then the stationary state  $p^\circ$  in which one type is absent and all other types play the “bad” strategy against all types must be in  $P$  (since the vector field vanishes on the straight line connecting this state with any other state in  $P$ ). But we know that if one type is absent in an inefficient stationary state, there is a tendency to move away from this state towards higher efficiency. (For instance, individuals of the absent type can “invade” the state  $p^\circ$  by playing the “good” strategy 1 with each other and strategy 2 with all other types.) So the set  $P$  is not asymptotically stable in any version of the co-ordination game with at least

two types.

An argument for efficiency in co-ordination games, very similar in spirit to the one given above, has recently been suggested in Kim and Sobel [10]. In that paper (of which we were unaware when we first wrote this paper), they apply a formal framework very similar to the one we set up in Section 3 to the study of pre-play communication in co-ordination games (or, in their terminology, common interest games). However, unlike us, they do not explicitly model the dynamics of an evolutionary process, but instead adopt a static evolutionary solution concept recently suggested by Swinkels [17], called *equilibrium evolutionary stability (EES)*. This concept differs from the more standard concept of evolutionarily stable strategies (ESS) in being a set-wise solution concept (like set-wise asymptotic stability) and in essentially requiring immunity only against invading populations which themselves are in (Nash) equilibrium (see Swinkels [17] for details). However, when one studies the inefficient outcome in co-ordination games, such invaders are easily produced since the potential invaders play the efficient Nash equilibrium strategy. Along these lines, one can show that the inefficient Nash equilibrium does not belong to an EES set. In fact, the argument is more or less exactly as above; if the inefficient equilibrium is to belong to an EES set then this set must also contain an equilibrium in which strategy 2 ("bad") is played but where there is an unsent message (the equivalent of an absent type in our framework), and this unsent message can be used to destabilize the inefficient outcome.

Kim and Sobel [10] are actually able to get an even stronger result. They show that *all* EES outcomes of a co-ordination game are Pareto efficient. The use of such set-valued stability concepts may thus provide powerful tools for eliminating inefficient equilibria. However, we feel that the embodied argument for elimination is not entirely convincing. After all, while it is true that the "bad" Nash equilibrium outcome in a co-ordination game (such as the one in Fig.3(a)) does not meet set-valued stability criteria, it is nevertheless both Lyapunov stable and has a nonempty basin of attraction. Hence, it is in fact a possible long-run outcome of the evolutionary process in some situations.<sup>7</sup>

---

<sup>7</sup>For example, let the initial population state be such that each type has only two (constant) behaviors, viz. play of either strategy 1 ("good") or 2 ("bad"), respectively, against all types. If the population share of each type playing "good" is sufficiently small, the population state will converge towards the stationary state in which all individuals

Kim and Sobel [10] provide a dynamic heuristic for the fact that the (static) EES-criterion eliminates the inefficient Nash equilibrium outcome along the following lines: there may be some “drift” within the EES set which leads away from this outcome towards a similar outcome with an un-sent message, but this latter outcome is (as argued above) unstable. However, such an idea of “drift” implicitly refers to some dynamics which is not explicitly modelled. One cause for such “drift” could be (small) costs associated with types *and/or* with the capacity to distinguish types. In fact, one can provide a cost-based story which, within a purely static evolutionary framework similar to that of Kim and Sobel [10], destabilizes the inefficient Nash equilibrium in the co-ordination game.

To see this, let, from now on, each individual be distinguished not only, as in Sections 3-5 above, by his type  $\tau$  and behavior  $\varphi$ , but also by his partitioning  $\pi$  of the set  $T$  of types. We formalize the notion that an individual cannot discriminate among types *within* cells of his partitioning by requiring that his behavior be constant on each cell of  $\pi$  (*i.e.*,  $\varphi$  has to be measurable with respect to  $\pi$ ). We now call a triplet consisting of a type  $\tau$ , a behavior  $\varphi$  and a partitioning  $\pi$ , a (*generalized*) *character*.

Furthermore, suppose each type  $\tau$  carries a different cost and that a partitioning with more cells costs more than a partitioning with fewer cells. We in fact only need (a) that there be at least *two* types which are more costly than the least costly type, and (b) that the (degenerate) partitioning with the full set  $T$  as its only cell costs less than all other partitionings (*i.e.*, that there is a cost associated with the capacity to discriminate at all). Both these requirements would seem to be generically met in any realistic setting.

Both kinds of costs, *i.e.*, the ones associated with being a type and with distinguishing types, respectively, are assumed to be lexicographically smaller than the payoffs of the underlying game. This allows us to employ the modified version of the ESS criterion, called MESS, suggested in a different evolutionary context in Binmore and Samuelson [5]. A strategy  $\sigma \in \Sigma$  in our expanded game  $\mathcal{G}$  is a *MESS* if, for any other strategy  $\mu \in \Sigma$ : either  $\sigma$  is a better reply than  $\mu$  to  $\sigma$ , or, if  $\mu$  is as good a reply against  $\sigma$  as  $\sigma$  is to itself, then  $\sigma$  is a better reply to  $\mu$  than  $\mu$  is to itself, or, if  $\mu$  is as good a reply to  $\sigma$  as  $\sigma$  is to itself *and*  $\sigma$  is as good a reply to  $\mu$  as  $\mu$  is to itself, then  $\sigma$  should not be more costly than  $\mu$ .

---

play “bad” against each other, just as in the standard evolutionary game theory setting.

Note that a MESS strategy  $\sigma \in \Sigma$  is a neutrally stable strategy (NSS) in the game  $\mathcal{G}$ , and therefore it is Lyapunov stable in the replicator dynamics on  $\Sigma$ .<sup>8</sup> We claim that in a co-ordination game of the kind discussed in Section 2 above there is *no* MESS in which only the “bad” strategy 2 is used. For if all types were to play the same strategy (2) against each other, then no (generalized) character which distinguishes between any two types could be present, since the state could then be invaded by a character which differs only by not making that distinction and hence is less costly. Therefore, all individuals in the population must have the whole set  $T$  as the only cell in their partitioning (just as in standard evolutionary game theory). But then carrying a type which is more costly than the least costly type is useless, since nobody makes any distinction in his behavior. Hence, all players have to be of the least costly type, which we may call type 1 without loss of generality (and we are completely back in the standard set-up of evolutionary game theory). However, if this were the case, then a new type could enter which would play strategy 2 against type 1 and strategy 1 against itself, and so earn more than all other players (by the lexicographic ordering of costs below payoffs). Hence, there is no MESS in which only the inefficient Nash equilibrium is played.

We finally note that with such costs as indicated above, there exists *no* MESS in games of the Prisoner’s Dilemma variety. In fact, this non-existence applies to any game which has a Nash equilibrium which is unique, pure, symmetric and Pareto dominated by another (possibly mixed) strategy combination. To see that there is no MESS in such a game, recall that a MESS is a NSS, and a NSS is a Lyapunov stable state, and hence every present type must play a Nash equilibrium with every present type, by Proposition 1. Hence, all individuals play the unique Nash strategy of the game  $G$ . By the same argument as given above concerning co-ordination games, only the least costly type, “type 1,” will be present, and no one will discriminate against any opponent. We now distinguish between two cases, one in which the dominating outcome involves some strategy  $x \in \Delta$  played against itself and one in which the dominating outcome involves two strategies,  $x'$  and  $x''$ , played against each other. In the first case, an absent type playing  $x$  against

---

<sup>8</sup>Neutral stability is the weakening of evolutionary stability that one obtains if the strict inequality in its definition is replaced by a weak inequality. Moreover, neutral stability implies Lyapunov stability in the replicator dynamics, see Weibull (1992).

itself and the Nash strategy against type 1 can invade. In the second case, two absent types, each playing the Nash equilibrium strategy against its own type and  $x'$  and  $x''$  against each other, respectively, can invade. In either case, there is no MESS.

Kim and Sobel [10] avoid this kind of instability by requiring that the invader should play an equilibrium strategy. This idea, due to Swinkels [17], is certainly an elegant way to restore existence. However, it seems unclear whether this is the most insightful route. We do feel that there is a real tendency towards efficiency in evolutionary games, and that this tendency sometimes conflicts with the Nash equilibrium logic, generating instances of non-existence. To recognize this conflict and to take to heart its implications may be more instructive than to rule out the possibility.

## Appendix

**Proposition:** *If a strategy  $i \in I$  is weakly dominated by some strategy  $\bar{x} \in \Delta$ , and  $u(\bar{x}, e^j)$  exceeds  $u(e^i, e^j)$ , then  $x_i(t)x_j(t) \rightarrow 0$  along any interior solution path to the replicator dynamics.*

*Proof:* Define  $w : \text{int}(\Delta) \rightarrow \mathfrak{R}$  by  $w(x) = -\ln(x_i) + \sum_{j=1}^m \bar{x}_j \ln(x_j)$ . Clearly  $w$  is continuously differentiable, and  $\dot{w}(x(t)) = u(\bar{x}, x(t)) - u(e^i, x(t)) \forall t$ . Since  $\bar{x}$  weakly dominates  $i$ , and  $x(t)$  is interior,  $w(x(t))$  increases monotonically over time. Suppose  $x_j(t)$  does not converge to zero. It suffices to show that  $w(x(t))$  then increases without bound, since the latter implies  $x_i(t) \rightarrow 0$ . Let  $\delta = u(\bar{x} - e^i, e^j)$ , a positive number. Let  $z = (x - x_j e^j)/(1 - x_j)$ , and note that  $z \in \Delta$  and  $x = (1 - x_j)z + x_j e^j$ . By bi-linearity of  $u$ ,  $\dot{w}(x(t)) = (1 - x_j)u(\bar{x} - e^i, z) + x_j u(\bar{x} - e^i, e^j)$ . Both these terms are non-negative, so

$$\limsup_{t \rightarrow \infty} \dot{w}(x(t)) \geq \delta \cdot \limsup_{t \rightarrow \infty} x_j(t)$$

*If  $w(x(t))$  were bounded, then the left hand side would be zero, a possibility which is excluded since  $x_j(t)$  does not converge to zero. Hence,  $w(x(t)) \rightarrow \infty$ .*

# Bibliography

- [1] Banerjee A. and J.W. Weibull, 1991, "Evolutionary selection and rational behavior", in Kirman A. and M. Salmon (eds.), **Rationality and Learning in Economics**, Blackwell (forthcoming).
- [2] Banerjee A. and J. Weibull, 1992, "Evolution and rationality; some recent game-theoretic results", mimeo., forthcoming in Allen B. (ed.), **The Proceedings of the Tenth World Congress of the International Economic Association**.
- [3] Bernheim D., 1984, "Rationalizable strategic behavior", **Econometrica** 52, 1007-1028.
- [4] Bhaskar V., 1991, "Noisy communication and the evolution of cooperation", mimeo., Delhi School of Economics.
- [5] Binmore K. and L. Samuelson, 1992, "Evolutionary Stability in Repeated Games Played by Finite Automata", **Journal of Economic Theory** 57, 278-305.
- [6] Bomze I., 1986, "Non-cooperative two-person games in biology: a classification", **International Journal of Game Theory** 15, 31-57.
- [7] Friedman D., 1991, "Evolutionary Games in Economics", **Econometrica** 59, 637-666.
- [8] Hirsch M. and S. Smale, 1974, **Differential Equations, Dynamical Systems, and Linear Algebra**, Academic Press.
- [9] Hofbauer J. and K. Sigmund, 1988, **The Theory of Evolution and Dynamical Systems**, London Mathematical Society Students Texts, Vol. 7, Cambridge University Press, Cambridge.

- [10] Kim Y.-G. and J. Sobel, 1991, "An evolutionary approach to pre-play communication", mimeo, University of Iowa and University of California at San Diego.
- [11] Nachbar J., 1990, "'Evolutionary' selection dynamics in games: convergence and limit properties", **International Journal of Game Theory** 19, 59-89.
- [12] Pearce D., 1984, "Rationalizable strategic behavior and the problem of perfection", **Econometrica** 52, 1029-1050.
- [13] Ritzberger K. and J. Weibull, 1992, "Evolutionary selection in normal-form games", mimeo., Stockholm University and Institute for Advanced Studies, Vienna.
- [14] Robson A.J., 1990, "Efficiency in evolutionary games: Darwin, Nash and the Secret Handshake", **Journal of Theoretical Biology** 144, 379-396.
- [15] Samuelson L., "Does evolution eliminate dominated strategies?", University of Wisconsin, mimeo.
- [16] Samuelson L. and J. Zhang, 1992, "Evolutionary Stability in Asymmetric games", **Journal of Economic Theory** 57, 363-391.
- [17] Swinkels J., 1992a, "Evolutionary Stability with Equilibrium Entrants", **Journal of Economic Theory** 57, 306-332.
- [18] Swinkels J., 1992b, "Adjustment dynamics and rational play in games", Stanford University.
- [19] Wärneryd K., 1991, "Evolutionary Stability in Unanimity Games with Cheap Talk", **Economics Letters** 36, 375-378.
- [20] Weibull J., 1992, "An introduction to evolutionary game theory", WP 347, The Industrial Institute for Economic and Social Research, Stockholm.