**Interim Report**    IR-09-089

# Public good games with incentives: The role of reputation

Karl Sigmund (ksigmund@iiasa.ac.at)
Hannelore De Silva (hannelore.desilva@wu.ac.at)

**Approved by**

Ulf Dieckmann
Leader, Evolution and Ecology Program

June 2010

# Public Good Games with Incentives: the role of reputation

Hannelore De Silva and Karl Sigmund

**Abstract**

Both the Trust Game and the Ultimatum Game reduce, in their most simplified versions, to a Public Good Game with an added incentive: namely a reward in the first case, and a sanction in the other. In this paper, the evolutionary game dynamics of these games is analyzed by means of the replicator equation. Positive and negative incentives have very different but complementary effects. We investigate the role of reputation, and show how occasional failures to contribute can lead to stabilizing cooperation.

## 1   A philosophical entente cordiale

In *Leviathan* (1651), the English philosopher Thomas Hobbes described life in the absence of a central authority as 'solitary, poore, nasty, brutish and short'. Selfish urges lead to 'such a war as is every man against every man.' The contemporary French philosopher Blaise Pascal held an equally dim view: 'Nous naissons injustes; car chacun tend à soi...La pente vers soi est le commencement de tout desordre en guerre, en police, en economie etc.' (We are born unfair; for everyone inclines towards himself...The tendency towards oneself is the origin of every disorder in war, polity, economy etc.)

In the following century, views on selfishness underwent a remarkable turn-about. The Scottish philosopher Adam Smith held that the selfish person works inadvertently for the public benefit. 'By pursuing his own interest he frequently promotes that of the society more effectually than when he really intends to promote it.' Greed promotes behavior beneficial to others. And most famously: 'It is not from the benevolence of the butcher, the brewer, or the baker, that we expect our dinner, but from their regard to their own self-interest. We address ourselves, not to their humanity but to their self-love, and never talk to them of our own necessities but of their advantages.'

An intriguingly similar view had been expressed, well before Smith, by Voltaire in his *Lettres philosophiques* (also known as *Lettres anglaises*): 'Il est bien vrai que Dieu aurait pu faire des creatures uniquement attentives

au bien d'autrui. Dans ce cas, les marchands auraient ete aux Indes par charite et le macon eut scie de la pierre pour faire plaisir à son prochain. Mais Dieu a etabli les choses autrement....C'est par nos besoins mutuels que nous sommes utiles au genre humain; c'est le fondement de tout commerce; c'est l'eternel lien des hommes.' ('Assuredly, God could have created beings uniquely interested in the welfare of others. In that case, traders would have been to India by charity, and the mason would saw stones to please his neighbor. But God designed things otherwise...It is through our mutual needs that we are useful to the human species; this is the grounding of every trade; it is the eternal link between men.')

The French term 'amour propre' certainly sounds a lot better than 'self-love'. Voltaire boldly claimed: 'Il est aussi impossible qu'une societe puisse se former et subsister sans amour propre, qu'il serait impossible de faire des enfants sans concupiscence, de songer à se nourrir sans appetit, etc. C'est l'amour de nous-même qui assiste l'amour des autres.' ('It is as impossible that a society could emerge and subsist without self-love than that people could produce children without lust, feed themselves without appetite, etc. Love for oneself assists the love for others.')

It is unknown whether Hobbes, during his time in Paris, ever met Pascal; but Smith most certainly had associated with Voltaire.

## 2  Public Goods and Private Incentives

So much for philosophical views on selfishness. They vary. But economic models make it clear that self-interested individuals will not act to achieve their group interest, except when prodded by incentives directed selectively towards individuals in the group, i.e. punishing exploiters or rewarding contributors (Olson 1965, see also Hardin 1966, Henrich and Boyd 2001, Sigmund 2007). Self-love is not always beneficial: it needs help to escape from the traps of social dilemmas. In this chapter, we investigate the role of reputation to promote an 'enlightened self-interest'. The importance of reputation as a kind of second (non-onetary) currency is well-established in economics literature, of course. Here we present a treatment based on evolutionary game dynamics (Hofbauer and Sigmund 1998, Nowak 2006). If players simply imitate what is successful in the long run, with nothing but self-interest in their mind, populations can evolve towards economically beneficially behavior.

We analyze a few basic models, starting with two scenarios which at first glance seem quite different, and which are well-known in behavioral game

theory as Trust Game and Ultimatum Game (Kagel and Roth 1995, Camerer 2003, Fehr and Camerer 2006).

Both are one-shot, two-person games. In both, a coin toss first decides who of the two players is the Proposer and who is the Responder. The Proposer is then endowed with a certain amount of money. In the Trust Game (Berg et al 1995), the Proposer can decide to donate part of this endowment to the Responder, knowing that it will be multiplied by a factor $r > 1$ by the experimenter. The Responder can then decide whether or not to return a part of this donation to the Proposer. This concludes the Trust Game. The Ultimatum Game (Güth et al 1982) does not take much time either. The endowment, in this case, is conditional. The Proposer has to offer a percentage $p$ of it to the Responder, and if the Responder accepts, the Proposer keeps the rest; but if the Responder declines, the experimenter withdraws the whole sum, so that both players gain nothing.

In the Trust game, a purely selfish Responder will never return anything, and a purely selfish Proposer, anticipating this, should offer nothing. In the Ultimatum game, a Responder's self-interest will accept any positive sum, since it is better than nothing. Accordingly, the Proposer should offer only a very small sum. In real experiments, the observed behavior differs considerably from these predictions of what a card-board 'homo economicus' ought to do. Indeed, in the Trust game, Responders often return a large part of their gift, and in the Ultimatum game, Responders often reject offers which they deem too small (Camerer 2003, Henrich 2006). Accordingly, Proposers in both types of games tend to transfer substantial proportions of their endowment, to both players' mutual benefit.

Both Trust and Ultimatum games are used to study norms of behavior, such as fairness and concern for one another. We shall study the evolutionary dynamics of simplified versions of these games, and then apply these results to address the issue of public goods with positive or negative incentives. Our main claim is that the concern for one's own reputation plays an essential role in causing us to deviate from what is prescribed for 'homo economicus', and hence to turn to economically more profitable behavior.

## 3   The Mini-Trust game

In a minimal variant of the Trust game, we assume that the Proposer has only to decide whether or not to donate a fixed amount $c$. Thus a Proposer has the choice between two moves $\mathbf{e}_1$ (donate) and $\mathbf{e}_2$ (defect). A Responder who receives a donation (i.e. the amount $b = rc$) has a choice between two

moves, namely to return a certain amount $\beta$ or not: these two moves will be denoted by $\mathbf{f}_1$ and $\mathbf{f}_2$. To make the game interesting, we will assume that $c < \beta < b$. In this case, if both players cooperate, both can make a gain. The payoff matrix is

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (\beta - c, b - \beta) & (-c, b) \\
\mathbf{e}_2 & (0, 0) & (0, 0)
\end{array}
\tag{1}
$$

Since the players are with equal probability in the role of Proposer and Responder, they are involved in a symmetric game. There exist four strategies, namely (a) the 'pro-social' strategy $G_1 = \mathbf{e}_1\mathbf{f}_1$ (donate, return); (b) the strategy $G_2 = \mathbf{e}_2\mathbf{f}_1$ (such a player does not donate, but returns a donation; (c) the asocial strategy $G_3 = \mathbf{e}_2\mathbf{f}_2$ (neither donate nor return); and finally (d) the strategy $G_4 = \mathbf{e}_1\mathbf{f}_2$ (such a player donates, but does not return). It is easy to compute the expected payoff values. But before doing this, we interpolate two brief sections on the replicator dynamics of two-role, two-strategy games (Hofbauer and Sigmund 1988, Sigmund et al 2001), in order to make this chapter self-contained.

## 4 The dynamics of two-role games

Let us consider a game with two roles I and II, and with two strategies for each role, which we denote by $\mathbf{e}_i$ and $\mathbf{f}_j$. The payoff matrix is

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (A, a) & (B, b) \\
\mathbf{e}_2 & (C, c) & (D, d)
\end{array}
\tag{2}
$$

Let us assume that a coin toss decides which role to assign to which player. The strategies for the resulting symmetric game will be denoted by $G_1 = \mathbf{e}_1\mathbf{f}_1$, $G_2 = \mathbf{e}_2\mathbf{f}_1$, $G_3 = \mathbf{e}_2\mathbf{f}_2$ and $G_4 = \mathbf{e}_1\mathbf{f}_2$. The payoff for a player using $G_i$ against a player using $G_j$ is given, up to the factor $1/2$ which we shall henceforth omit, by the $(i, j)$-entry of the matrix

$$
M = \begin{pmatrix}
A + a & A + c & B + c & B + a \\
C + a & C + c & D + c & D + a \\
C + b & C + d & D + d & D + b \\
A + b & A + d & B + d & B + b
\end{pmatrix}.
\tag{3}
$$

Let us assume that players tend to imitate successful individuals, and hence occasionally switch from one strategy to another. They compare their average payoff with that of another player and adopt that player's strategy with a

probability proportional to the payoff difference, if it is positive (if not, they do not switch). Since the payoffs depend on the state of the (well-mixed) population, given by the frequencies $x_i(t)$ of the strategies $G_i$, this yields an evolutionary dynamics in the state space $S_4 = \{(x_1, x_2, x_3, x_4) \in R_+^4 : x_1 + ... + x_4 = 1\}$. It is given by the replicator equation

$$\dot{x}_1 = x_i[(M\mathbf{x})_i - \bar{M}], \tag{4}$$

where $\bar{M} = x_1(M\mathbf{x})_1 + ... + x_4(M\mathbf{x})_4$ is the average payoff in the population. Since the dynamics are unaffected if one modifies the payoff matrix $M$ by replacing $m_{ij}$ by $m_{ij} - m_{1j}$, we can use the matrix

$$\begin{pmatrix} 0 & 0 & 0 & 0 \\ R & R & S & S \\ R+r & R+s & S+s & S+r \\ r & s & s & r \end{pmatrix}. \tag{5}$$

with $R := C - A$, $r := b - a$, $S := D - B$ and $s := d - c$.

# 5   Staying in the saddle

We shall denote matrix (5) again by $M$. It has the property that

$$m_{1j} + m_{3j} = m_{2j} + m_{4j} \tag{6}$$

for $j = 1, 2, 3, 4$. Hence

$$(M\mathbf{x})_1 + (M\mathbf{x})_3 = (M\mathbf{x})_2 + (M\mathbf{x})_4 \tag{7}$$

holds for all $\mathbf{x}$. From this follows easily that the function $V = x_1 x_3 / x_2 x_4$ satisfies

$$\dot{V} = V[(M\mathbf{x})_1 + (M\mathbf{x})_3 - (M\mathbf{x})_2 - (M\mathbf{x})_4] = 0 \tag{8}$$

in the interior of $S_4$, and hence the value of $V$ remains unchanged along every orbit.

Hence the interior of the state simplex $S_4$ is foliated by the surfaces

$$W_K := \{\mathbf{x} \in S_4 : x_1 x_3 = K x_2 x_4\}, \tag{9}$$

with $0 < K < \infty$. These are saddle-like surfaces which are spanned by the quadrangle of edges $G_1 G_2$, $G_2 G_3$, $G_3 G_4$ and $G_4 G_1$ joining the vertices of the simplex $S_4$ (see Fig 1).

5

Figure 1: The state space $S_4$ (a simplex with four corners $G_i$, i=1,2,3,4, corresponding to the four strategies of a symmetrized two-roles, two-strategies game), and a saddle-like surface $W_K$ spanned by the edges $G_1 \to G_2 \to G_3 \to G_4 \to G_1$ (see text). The evolving states remain on their initial surface $W_K$. If there exist fixed points in the interior of the state space, they form a line intersecting each $W_K$.

The orientation of the flow on the edges can easily be obtained from the previous matrix. For instance, if $R = 0$, then the edge $G_1 G_2$ consists of fixed points. If $R > 0$, the flow along the edge points from $G_1$ towards $G_2$ (in the absence of the strategies $G_3$ and $G_4$, the strategy $G_2$ dominates $G_1$), and conversely, if $R < 0$, the flow points from $G_2$ to $G_1$.

Generically, the parameters $R, S, r$ and $s$ are non-zero. This corresponds to 16 orientations of the quadrangle $G_1 G_2 G_3 G_4$, which by symmetry can be reduced to 4 (see Fig 2). Fixed points in the interior of the simplex $S_4$ must satisfy $(M\mathbf{x})_i = 0$ for $i = 2, 3, 4$ (since $(M\mathbf{x})_1$ trivially vanishes). This implies for $S \neq R$

$$x_1 + x_2 = \frac{S}{S - R}, \tag{10}$$

and for $s \neq r$

$$x_1 + x_4 = \frac{s}{s - r}. \tag{11}$$

Such solutions lie in the simplex if and only if $RS < 0$ and $rs < 0$, which corresponds to the orientations (c) and (d) of the quadrangle spanning the saddle-like surfaces $W_K$. If this is the case, one obtains a line of fixed points which intersects each $W_K$ in exactly one point (see Fig. 1). The solutions can be written as

$$x_i = m_i + \xi \tag{12}$$

6

for $i = 1, 3$ and

$$x_i = m_i - \xi \tag{13}$$

for $i = 2, 4$, with $\xi$ as parameter and

$$\mathbf{m} = \frac{1}{(S - R)(s - r)}(Ss, -Sr, -Rr, Rs) \in W_1. \tag{14}$$

Figure 2: The four generic orientations of the quadrangles spanning the saddle-like surfaces. The orientations depend on the signs of $R, S, r$ and $s$ (see text). In cases (c) and (d), there exists a fixed point in the interior of $W_K$.

**Figure 2 should be
approximately here.**

# 6 Farewell to Trust

For the corresponding payoff matrix, we obtain $R = c - \beta < 0$, $r = \beta > 0$, $S = c > 0$ and $s = 0$ (see Fig 3). If $x_3 = x_4 = 0$, i.e., if everyone in the population is ready to return a donation, it is best to donate, i.e, $G_1$ dominates $G_2$. If $x_2 = x_3 = 0$, i.e., if donations can be taken for granted, then it is best not to return it, i.e., $G_4$ dominates $G_1$. If $x_1 = x_2 = 0$, i.e., if no one ever returns a donation, then $G_3$ dominates $G_4$, i.e., it is best not to donate. Finally, if $x_1 = x_4 = 0$, i.e., if nobody ever donates, then it does not matter whether one is willing to return a donation or not. In this case, the state of the population is a fixed point. Neither $G_2$ nor $G_3$ has an advantage.

It is easy to see that the segment $QG_3$, where

$$Q = (0, \frac{c}{\beta}, \frac{\beta - c}{\beta}, 0), \tag{15}$$

**Figure 3 should be
approximately here.**

Figure 3: Dynamics on a saddle-like surface for the Trust game (or for a Public Good game with reward). The edge $G_1G_4$ consists of fixed points, the segment $G_1Q$ of stable fixed points which are Nash equilibria.

consists of saturated fixed points, i.e., of Nash equilibria. Indeed, for $x_1 = x_4 = 0$, both $(M\mathbf{x})_1$ (which is normalized to 0) and $(M\mathbf{x})_4$ are smaller than the average payoff $\bar{M} = (M\mathbf{x})_2 = (M\mathbf{x})_3 = c - \beta x_2$. The flow along the edges leads from $G_2$ to $G_1$, from there to $G_4$, and then to $G_3$. All orbits in the interior converge to the segment $QG_3$ for $t \to +\infty$ and to the segment $QG_2$ for $t \to -\infty$. Thus the population will, in the long run, consist only of players who, as Proposers, never donate (and consequently, as Responders, never return anything). From the economic viewpoint, the minimal version of the Trust game does not take off: no donations, no paybacks.

# 7 Ultimate offers

We now turn to the Ultimatum game. It is simple enough, but we shall simplify it even further (cf. Nowak, Page and Sigmund 2000), and assume that the Proposer has only a choice between offering a high percentage $h$ (for instance, 45 percent) or a low percentage $l$ (for instance 15 percent), with $0 < l < h < 1$. The Responder could, in principle, accept both offers, one of them, or none. Again, we simplify by assuming that he has to choose between two strategies only: the strategy denoted by $h$, which consists in accepting the high offer only, or the strategy denoted by $l$, which consists in accepting both possible offers.

In this reduced version of the Ultimatum game, the two strategies for role I, namely $\mathbf{e}_1$ and $\mathbf{e}_2$, are given by the offers $h$ and $l$; and the two strategies $\mathbf{f}_1$ and $\mathbf{f}_2$ for role II will again denoted by $h$ and $l$, for convenience; these

strategies correspond now the Responder's aspiration levels. The payoff matrix is given by

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (1-h, h) & (1-h, h) \\
\mathbf{e}_2 & (0, 0) & (1-l, l)
\end{array}
\tag{16}
$$

The strategy $G_1$ corresponds to $(h, h)$: high offers, and a high aspiration level. We may term it as the *fair* strategy. By contrast, $G_3 = (l, l)$ epitomizes the selfish strategy. It leads to the acceptance of any positive offer, and aims to part with as little as possible. The strategy $G_2 = (l, h)$ is paradoxical: it offers little, but insists on a high offer. $G_4$, finally, makes a good offer, but accepts a low offer. For want of a better term, we call it the *mild* strategy. The payoff parameters are $R = h - 1 < 0$, $r = 0$, $S = h - l > 0$ and $s = l > 0$. The selfish strategy is dominated by the mild strategy, which is dominated by the paradoxical strategy, which in turn is dominated by the fair strategy; but the mild and the fair strategies are equivalent, in the absence of the other two strategies, one does as well as the other: all offers are fair, and the average payoff is $1/2$.

There exist no fixed points in the interior of $S_4$. Indeed, whenever $x_2 > 0$ or $x_3 > 0$, we have $(M\mathbf{x})_4 > (M\mathbf{x})_1$ and hence both ratios $x_4/x_1$ and $x_3/x_2$ are increasing. On each surface $W_K$, the flow is as shown in Fig 4. On the edge $x_2 = x_3 = 0$, all points are fixed points. If $x_1 < \frac{h-l}{1-l}$, then both $(M\mathbf{x})_2$ and $(M\mathbf{x})_3$ are larger than $\bar{M}$. Let us denote by $\mathbf{Q}$ the point $(\frac{h-l}{1-l}, 0, 0, \frac{1-h}{1-l})$. Then the symmetric Nash equilibria of the game are those on the segment $G_3Q$, and the vertex $G_1$. We note that on the edge $x_2 = x_4 = 0$, there exists another fixed point $P$, with coordinates $(h, 0, 1 - h, 0)$. In a population with selfish and fair players only, we have a bistable competition. The fair strategy is risk-dominant (i.e., a population consisting in equal numbers of selfish and fair players will see fair players win) if $h < 1/2$.

The orbits in the interior of $S_4$ either converge to $G_3$, or else to the set of Nash equilibria, as shown in Fig 4. If we assume that random shocks occasionally perturb the state of the population, we will expect that they lead to neutral drift along the edge $x_2 = x_3 = 0$. As soon as $x_1 < \frac{h-l}{1-l}$, a random perturbation sending the state into $int S_4$ will cause the fixation of $G_3$. This implies that eventually, the population consists of selfish players only. Thus evolutionary game theory leads to the same prediction as classical game theory; both are in contrast to experimental evidence.

9

**Figure 4 should be
approximately here.**

Figure 4: Dynamics on a saddle-like surface for the Ultimatum game (and Public Good game with punishment). The edge $G_1 G_4$ consists of fixed points, the segment $G_1 Q$ of stable fixed points which are Nash equilibria.

# 8 Bifurcation through Reputation

So far, we have considered conditions of strict anonymity. Let us now assume that with some (possibly small) probability, players may know their co-player by reputation, and in particular may know about the offers previously accepted by that co-player. Let us furthermore assume that occasionally, players offer less than they usually would, if they have reason to believe that they can get away with it; more precisely, if they know that their co-player has previously accepted low offers. The two assumptions seem reasonable enough: they only require some information about other players in the group, and a touch of opportunistic selfishness. In that case, accepting a low offer can have the regrettable consequence that one is offered less, in future games.

In order to analyze this situation, let us assume that $\mu > 0$ is the probability that a 'fair' $(h, h)$ Proposer encountering a mild $(h, l)$ Responder knows that this player accepts a low offer, and consequently offers $l$ instead of $h$. This yields the payoff matrix

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (1-h, h) & (1-h+\mu(h-l), h-\mu(h-l)) \\
\mathbf{e}_2 & (0, 0) & (1-l, l)
\end{array}
\tag{17}
$$

which differs from (1) in one position only, by the term $\mu(h - l)$ which can be arbitrarily small. It can be viewed as a perturbation of the previous game, due to the effect of reputation. The corresponding symmetrized game

10

(5) is now given by $R = h - 1$, $r = -\mu(h - l)$, $S = (h - l)(1 - \mu)$ and $s = l$. For $\mu < 1$, we have $R < 0$, $S > 0$, $s > 0$ (as before) and $r > 0$ (while we had $r = 0$ in the unperturbed case). This yields now a generic case, corresponding to case (c) in Fig 2. There exists a line of fixed points in the interior of the state space $S_4$. Each of the surfaces $W_K$ (for $K > 0$) intersects this line in a saddle point. For $\mu \to 0$, the point $\mathbf{m}$, and with it all interior fixed points, converge to the point $Q$ on the edge $G_1G_4$. The dynamics on each surface $W_K$ is bistable, the vertices $\mathbf{e}_1$ and $\mathbf{e}_3$ are the attractors (see Fig 5). Hence, depending on the initial condition, the population will either converge to the fair or to the selfish strategy.

# 9  Public Goods with Punishment

In a simple form of the Public Goods game, each of the $N$ players participating in the game has the possibility of contributing a fixed amount $c$ to the common pool. The experimenter multiplies each player's contribution by a factor $r > 1$, and divides the resulting amount equally among all other $N - 1$ players participating in the game.

For $N = 2$, this is a Prisoner's Dilemma game game: both players can decide whether or not to send a donation $b = rc$ to the other player, at a cost $c$ for themselves. The dominant solution is to defect. But let us now introduce a second stage to this game, by allowing the players to punish defectors. We shall assume that the sanction consists in imposing a fine of size $\beta$. This fine is not collected by the punishing player. On the contrary, the punisher has to pay a fee, which costs him an amount $\gamma$. The first stage of the game offers scope for altruism (helping another player at a cost for oneself), and the second stage scope for spite (harming the other player at a cost for oneself). Obviously, in both stages, the dominating solution is to avoid the cost. A selfish player should defect in the first stage, and refuse to punish in the second stage.

If we assume that players can impose their fine conditionally, fining only those who have failed to help them, the long-term outcome will be, as before, that no pro-social behavior emerges (see Sigmund, Hauert and Nowak 2001). Indeed, let us label with $\mathbf{e}_1$ those players who cooperate by sending a donation to their co-player, and with $\mathbf{e}_2$ those who do not, i.e. who defect; similarly, let $\mathbf{f}_1$ denote those who punish defectors, and $\mathbf{f}_2$ those who do not.

**Figure 5 should be
approximately here.**

Figure 5: Dynamics on a saddle-like surface for an Ultimatum game with reputation (or for a Public Good game with reputation). The dynamics is bistable, the pro-social state $G_1$ and the asocial state $G_3$ are attractors.

The payoff matrix is given by

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (-c, b) & (-c, b) \\
\mathbf{e}_2 & (-\beta, -\gamma) & (0, 0)
\end{array}
\tag{18}
$$

Here, the first number in each entry is the payoff for the corresponding row player, and the second number for the column player. We have used the same notation, as for two-role games, although the situation is completely symmetric: instead of being either in one role or in the other, a player is first in one role and then in the other. Despite this difference, we can apply the same method as before. Indeed, each strategy for this two-stage game must specify what to do in the first stage, and what to do in the second. Hence, it is given by a pair $\mathbf{e}_i \mathbf{f}_j$ (with $i, j \in \{1, 2\}$). As in section 3, we denote the resulting four strategies with $G_1 = \mathbf{e}_1 \mathbf{f}_1$, $G_2 = \mathbf{e}_2 \mathbf{f}_1$, $G_3 = \mathbf{e}_2 \mathbf{f}_2$ and $G_4 = \mathbf{e}_1 \mathbf{f}_2$. The strategy $G_1$ corresponds to the 'pro-social' behavior: to give help, and to punish those who don't. $G_3$ is the selfish strategy which avoids any costs: a player using it does not help the co-player, and expects no help. $G_2$ can again be viewed as paradoxical: a $G_2$-player defects, but punishes a co-player who defects. Finally, $G_4$ can again be viewed as a 'mild' strategy: a $G_4$ player sends a donation to the co-player but does not react if this is not reciprocated.

12

# 10  Dynamics with reputation

We can follow the same approach as before, and obtain $R = c - \beta$, $S = c$, $r = 0$ and $s = \gamma$. Again, the manifolds $W_K = \{\mathbf{x} \in S_4 : x_1 x_3 = K x_2 x_4$ are invariant (for $K > 0$) and the dynamics is as in Fig 3. In fact, the Ultimatum mini-game can be viewed as a special case, with $\gamma = l$, $\beta = 1 - l$, and $b = c = h - l$. Intuitively, this simply means that in the Ultimatum game, the donation consists of making the high offer instead of the low offer. The benefit to the recipient (i.e. the Responder) $h - l$ is equal to the cost to the donor (i.e. the Proposer). The punishment consists of refusing the offer. This costs the Responder $l$ (the amount offered) and punishes the Proposer by the amount $1 - l$, which is large if the offer is low.

The fixed points in $W_K$ are the corners $\mathbf{G}_i$ and the points on the edge $\mathbf{G}_1 \mathbf{G}_4$. $\mathbf{G}_3$ is a Nash equilibrium, $\mathbf{G}_2$ is not. A point $\mathbf{x}$ on the edge $\mathbf{G}_1 \mathbf{G}_4$ is a Nash equilibrium whenever $x_1 \geq c/\beta$. Thus if $c > \beta$, $\mathbf{G}_3$ is the only Nash equilibrium. This case is of little interest. From now on, we restrict our attention to the case $c < \beta$: the fine costs more than the donation. We denote the point $(c/\beta, 0, 0, (\beta - c)/\beta)$ with $\mathbf{Q}$ and see that the closed segment $\mathbf{Q}\mathbf{G}_1$ consists of Nash equilibria. In the long run, in spite of the segment of Nash equilibria, random shocks will ultimately establish the asocial state $\mathbf{G}_3$.

Still following the parallel with the Ultimatum game, let us assume that with a probability $\mu$, cooperators (i.e. $\mathbf{e}_1$-players) defect against non-punishers, i.e. $\mathbf{f}_2$-players. (Hence $\mu$ is the probability that (1) the $\mathbf{f}_2$-type becomes known and (2) the $\mathbf{e}_1$-type decides to defect). The payoff matrix becomes

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (-c, b) & (-c(1 - \mu), b(1 - \mu)) \\
\mathbf{e}_2 & (-\beta, -\gamma) & (0, 0)
\end{array}
\tag{19}
$$

We obtain $R = (c - \beta) < 0$, $S = c(1 - \mu) > 0$, $s = \gamma > 0$ and $r = -b\mu < 0$. Thus the edge $\mathbf{G}_1 \mathbf{G}_4$ consists no longer of fixed points, but of an orbit converging to $\mathbf{G}_1$. The dynamics is as in Fig 5. On each saddle-like surface $W_K$, and therefore in the whole interior of the state space $S_4$, the dynamics is bistable, with attractors $\mathbf{G}_1$ and $\mathbf{G}_3$. Depending on the initial condition, every orbit converges to one of these two attractors, namely the asocial state $\mathbf{G}_3$ (no contributions, no punishment) and the pro-social regime $\mathbf{G}_1$ (cooperate, punish defectors).

# 11 Revealing errors

The previous model is, in a certain sense, incomplete. Indeed, it essentially depends on altering the dynamics on the edge $\mathbf{G_1}\mathbf{G_4}$ by introducing the reputation effect. But on that edge, the population consists of two types only, both contributing to the public good. How should players learn whether the co-player is of type $\mathbf{f}_1$ or $\mathbf{f}_2$, i.e. willing to punish a defector, or not? Even if each player plays many rounds of the game, no defection ever arises.

There are several ways to deal with this question. One possibility would be to assume that players learn about their co-players' propensity to punish from other sources. It seems not unlikely that we can get a good idea about the irascibility or meekness of our co-players by watching their interactions with noisy children or their reactions to the daily news, rather than merely from observing how they act in the donation game. But it is probably better to complete the model without appealing to other interactions.

The simplest approach is to introduce the possibility of errors. Let us assume that player play the game repeatedly, and that players intending to donate will, with a certain probability $\epsilon$, fail to implement their intention. (This could be due to a mistake, or to a lack of resources.) In the absence of reputation, this yields the following payoff structure:

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (-(1-\epsilon)c - \epsilon\beta, (1-\epsilon)b - \epsilon\gamma) & (-(1-\epsilon)c, (1-\epsilon)b) \\
\mathbf{e}_2 & (-\beta, -\gamma) & (0,0)
\end{array}
\tag{20}
$$

Compared with the situation in the previous section, $s$ remains unchanged, whereas $R$ and $S$ are multiplied by $(1 - \epsilon)$, which does not affect the sign, and hence conserves the dynamics on the corresponding edge. But $r$ is now equal to $\epsilon\gamma$, and hence positive. This means that on the edge $\mathbf{G_1}\mathbf{G_4}$, the flow points towards $\mathbf{G_4}$: punishment is dominated. As a result, we obtain a dynamics as in case (b) of Fig 2. All orbits in the interior of the simplex $S_4$ converge to the vertex $\mathbf{G_3}$. The asocial type wins.

Now let us introduce reputation. For simplicity, we will assume that players who know that their co-player is not of the punishing type never donate. (It would suffice to assume that they defect with a small probability). The parameter $\mu$, then, is simply the probability to learn that the co-player has, once in the past, failed to punish a defector. If we assume perfect information, this reduces to the probability that the co-player has encountered a defection. On the edge $x_2 = x_3 = 0$, all players are willing to donate, and a defection occurs only by mistake. The probability that the co-player, in his $k$ previous rounds, never faced a mistaken defection is $(1 - \epsilon)^k$. If the num-

ber of rounds is distributed geometrically, with a constant probability $w < 1$ for a further round, then $w^k(1 - w)$ is the probability that the co-player has experienced $k$ rounds. This means that

$$\mu = \frac{w\epsilon}{1 - w(1 - \epsilon)}. \tag{21}$$

If we assume that a player does not donate if he knows that he can get away with it (or if he commits an error), this yields

| | $\mathbf{f_1}$ | $\mathbf{f_2}$ |
|---|---|---|
| $\mathbf{e_1}$ | $(-(1 - \epsilon)c - \epsilon\beta, (1 - \epsilon)b - \epsilon\gamma)$ | $(-(1 - \epsilon)(1 - \mu)c, (1 - \epsilon)(1 - \mu)b)$ |
| $\mathbf{e_2}$ | $(-\beta, -\gamma)$ | $(0, 0)$ |

$$\tag{22}$$

We see that $r = \epsilon\gamma - \mu(1 - \epsilon)b$ is negative if

$$\gamma < \frac{w(1 - \epsilon)b}{1 - w(1 - \epsilon)}, \tag{23}$$

i.e., if the fee for punishing the defector is not too high.

Of course this can also be applied to the Ultimatum game. In that case, $r = \epsilon\gamma - \mu(1 - \epsilon)b$ is negative if

$$l < w(1 - \epsilon)h, \tag{24}$$

i.e. if the low offer is sufficiently smaller than the high offer.

## 12   Public Goods with rewards

Let us now consider a public good game (still with $N = 2$ players only), but assume that the players have, in a second phase of the game, the option of rewarding contributors. Thus we consider a positive rather than a negative incentive. We shall assume that players who reward their donors have to pay a cost $\gamma$, and that the rewarded player receives an amount $\beta$ (if $\beta = \gamma$ this is simply a payback). We assume $0 < c < \beta$ and $0 < \gamma < b$. If $\mathbf{e_1}$ and $\mathbf{e_2}$ are the two options for the first stage (to contribute or not), and $\mathbf{f_1}$ and $\mathbf{f_2}$ for the second stage (to reward donors or not), then the payoff structure is given by

| | $\mathbf{f_1}$ | $\mathbf{f_2}$ |
|---|---|---|
| $\mathbf{e_1}$ | $(\beta - c, b - \gamma)$ | $(-c, b)$ |
| $\mathbf{e_2}$ | $(0, 0)$ | $(0, 0)$ |

$$\tag{25}$$

The minimal variant of the Trust game, introduced in section 3, can be viewed as a special case (making the usual analogy between a two-role game

and a two-stages game). There exist four strategies, namely (a) the 'pro-social' strategy $G_1 = \mathbf{e}_1\mathbf{f}_1$ (donate, reward); (b) the strategy $G_2 = \mathbf{e}_2\mathbf{f}_1$ (such a player does not donate, but rewards a donor); (c) the asocial strategy $G_3 = \mathbf{e}_2\mathbf{f}_2$ (which neither donates nor rewards); and finally (d) the strategy $G_4 = \mathbf{e}_1\mathbf{f}_2$ (such a player donates, but does not reward). For the corresponding payoff matrix (5), we obtain $R = c - \beta < 0$, $r = \gamma > 0$, $S = c > 0$ and $s = 0$ (see Fig 3). The outcome is exactly the same as for the trust game. Thus the population will, in the long run, consist only of players who always defect (and consequently never reward).

Let us now introduce reputation effects into the Public Goods game with rewards. We shall assume that with a small likelihood $\mu$, cooperators defect if they know that the other player is not going to reward them, i.e. is of type $\mathbf{f}_2$. ($\mu$ is the probability that (1) the $\mathbf{f}_2$-type becomes known and (2) the $\mathbf{e}_1$-type decides to defect). Similarly, we denote by $\nu$ the likelihood that defectors cooperate if they know that they will be rewarded. ($\nu$ is the probability that (1) the $\mathbf{f}_1$-type becomes known and (2) the $\mathbf{e}_2$-type reacts and donates). This yields the payoff matrix

$$
\begin{array}{c|cc}
 & \mathbf{f}_1 & \mathbf{f}_2 \\
\hline
\mathbf{e}_1 & (\beta - c, b - \gamma) & (-c(1-\mu), b(1-\mu)) \\
\mathbf{e}_2 & ((\beta - c)\nu, (b - \gamma)\nu) & (0,0)
\end{array}
\tag{26}
$$

Now $R = (c-\beta)(1-\nu) < 0$, $S = c(1-\mu) > 0$, $r = \gamma - b\mu$ which is positive if $\mu$ is small, and $s = (\gamma - b)\nu$, which is negative. It is this last condition that differs from the unperturbed system. The edge $G_2G_3$ no longer consists of fixed points. Instead, $G_3$ is dominated by $G_2$: if players can acquire a reputation for rewarding donations, this can motivate co-players to donate. The essential parameter, therefore, is $\nu$.

Let us begin by assuming that $\mu$ is small, so that $r$ is positive. For $\nu > 0$, the flow on the edge $\mathbf{G}_2\mathbf{G}_3$ leads towards $\mathbf{G}_3$, so that the frame spanning the saddle-like surfaces $W_K$ is cyclically oriented (see Fig 6). As before, there exists now a line of fixed points in the interior of $S_4$. On each saddle-like surface $W_K$, the orbits rotate around this fixed point; they spiral towards it for $0 < K < 1$ and away from it for $K > 1$. The surface $W_1$ consists of periodic orbits.

We stress the highly unpredictable dynamics if $\nu > 0$ and $\mu$ small. For one half of the initial conditions, the replicator dynamics sends the state towards the line of fixed points. But there, random fluctuations will eventually lead to the other half of the simplex, where the replicator dynamics leads to the heteroclinic cycle $\mathbf{G}_1\mathbf{G}_4\mathbf{G}_3\mathbf{G}_2$. The population seems glued for a long time to one strategy, then suddenly switches to the next, remains there for

a still longer time etc... However, an arbitrarily small random shock will send the state back into the half-simplex where the state converges again to the line of fixed points, etc. Not even the time averages of the frequencies of strategies converge. One can only say that the most probable state of the population is either monomorphic (i.e. close to one corner of $S_4$) or else close to the attracting part of the line of fixed points (all four types present, the proportion of cooperators larger among rewarders than among non-rewarders, and – if the values $\nu$ is small – a frequency of rewarders close to $c/\beta$, and a frequency of donors which is small).

**Figure 6 should be
approximately here.**

Figure 6: Dynamics on a saddle-like surface for the Trust game with reputation ($\mu$ small, $\nu > 0$). The edges are cyclically oriented. For $W_1$, the orbits are periodic. The orbits on $W_K$ converge either to the inner fixed point or to the boundary, depending on whether $0 < K < 1$ or $K > 1$.

Let us note that we encounter the same problem as for the Public Good with Punishment. If $x_1 = x_4 = 0$, then nobody ever donates. In this case, how should the $\mathbf{f}_1$-trait (rewarding donors) ever reveal itself? The assumption that occasionally players commit errors is far less plausible as in the previous case, since donating inadvertently is far less likely than failing in the intention to donate.

Finally, let us briefly consider the case when the fact that a player does not reward has a high probability to become publicly known. In that case, it is unlikely that such a player receives a donation. This means that $\mu$ is close to 1, and hence that the parameter $r = \gamma - b\mu$ is negative. In that case, all orbits in the saddle-like surface $W_K$ converge to $G_1$ (see Fig. 7): the social

17

strategy wins.

Figure 7: Dynamics on the saddle-like surface for the Trust game if $\mu > \beta/b$. In that case, all orbits converge to $G_1$.

**Figure 7 should be
approximately here.**

# 13   Larger groups

So far we have only considered games with two players. Many economic interactions, and in particular many joint enterprises, involve more than two actors. In section 9 we have introduced a so-called *others only* version of the Public Good game with $N$ players. Each player's contribution is multiplied by $r > 1$ and divided equally among all $N-1$ other players. In another version, we can assume that it is divided among all $N$ players, so that if a player contributes $c$, then a part $\frac{r}{N}$ is returned to the donor. In the simplest case, when each contribution is of the same value $c$ and if $N_c$ players contribute, then the total amount $cN_c$ is multiplied by $r > 1$ and divided equally among all $N$ participants. A social dilemma holds if $r < N$. In alternative models, the total amount is a non-linear function of the number of contributors.

Similarly, there are many ways of modeling punishment. In the simplest approach, each punisher pays a fee $\gamma$ to inflict a fine $\beta$ upon each defector. The resulting game dynamics is like that with two players (Hauert et al 2004). If random shocks occasionally perturb the system, then in the long run, the asocial strategy (no contribution, no punishment) dominates the population. Again, the situation can be redressed if we assume that players can obtain information about the type of their co-players, and that contributors occasionally yield to the temptation of exploiting their co-players if they know that they can get away with it (i.e., that there are few or no punishers in their group).

18

With positive rewards, the situation is again similar to that of a two-person game, at least for a large set of paramter values.

# 14   Discussion

It is unlikely that one-shot interactions between anonymous players, such as the Ultimatum game or the Trust game, play a prominent role in human economy. In fact, their artificiality is an advantage for experiments. From early on, most experiments in physics or physiology are similarly based on artificial situations, such as a feather in a vacuum tube etc.

On the other hand, some everyday parallels to Trust and Ultimatum games exist. For instance, sellers who fix a (non-negotiable) price tag to an object displayed in their shopwindow are proposing an offer to the passersby, who can reject it or not. This has similarities with the Ultimatum game. And individuals entrusting their banker with money are engaging in a transaction similar to a Trust game. In everyday life, we often see that contributions to the public good are encouraged by heavily fining free-riders, etc (Henrich 2006, Ostrom and Walker 2003). On the other hand, there are essential differences between the games and the real-life parallels. For instance, many passersby will look into the shopwindow, whereas the Ultimatum game has only one Responder (if there are several, the outcome is drastically altered).

In each of these games, reputation can play an essential role in boosting the economically advantageous strategy (just as in indirect reciprocity, see (Nowak and Sigmund 2005, Wedekind and Milinski 2000). Reputation requires an information flow in the population. This information flow extends the knowledge obtained through the games that are personally experienced by a player, and usually relies on gossip. For instance, we have seen in section 12 that as soon as it is safe to assume that a funds manager who returns less than the investment becomes publicly known, the social strategy (for the clients, to invest, and for the managers, to return more than that investment to the clients) is a global attractor. Another example concerns internet trading, such as e-Bay. It relies heavily on the possibility that clients can rate their former partners. Another argument stems from psychology. If individuals feel unfairly treated, they often vent their emotions to a large audience (see e.g. Xiao and Houser 2005). Anger is loud. The logic behind this is clear: rejecting an unfair but positive offer involves costs, which can only be recouped if they prevent others from making unfair offers. If you take the trouble of getting emotional, you should make it known.

The importance of information has been displayed in a neat experiment

based on two treatments of the Ultimatum game (Fehr and Fischbacher 2003). Each player engages in several such games (always against a different partner, of course). In one treatment, players are anonymous. In the other treatment, players have pseudonyms and know that their decisions, as Responder, will be made known to their future Proposers. Aspiration levels are significantly higher in the second treatment. Players are more likely to reject offers. They seem to expect that if they once accept a low offer, they run a high risk of encountering such offers again and again. (See also http://homepage.univie.ac.at/hannelore.brandt/ultimatum/ for on-line computer simulations, cf. Figs 8 a and b).

Figure 8: Two variants of individual-based simulations on the Ultimatum game. In both cases, 1000 fictitious players with randomly chosen strategies $(p, q)$ (where $p$ is the size of the offer and $q$ the aspiration level) each play 50 games against randomly chosen co-players. Then, the frequencies of the strategies are updated according to the replicator dynamics. This is repeated for many 'generations'. Left: players are anonymous. The population average of the $(p, q)$-values starts out close to the center ($p$ and $q$ close to 50 percent). After a drop in the $p$-value, the population average converges back to the diagonal and then inches along the diagonal towards $(0, 0)$. Right: players know the past of their co-players, and offer the minimum of their own $p$-value and the offers previously accepted by their co-player. The evolution begins similarly. But then, when the population average has returned to the diagonal, the $p$ and $q$-values creep up, not down, and reach a value slightly below 50 percent.

**Figure 8 should be
approximately here.**

Of course, even if players know perfectly well that their action is not observed, they often act as if it were. The lingering suspicion that despite double-blind conditions etc. someone could be watching, is neatly captured in a series of experimental papers that show that the mere picture of an eye (on a poster, or on a computer screen) can activate a subconscious concern
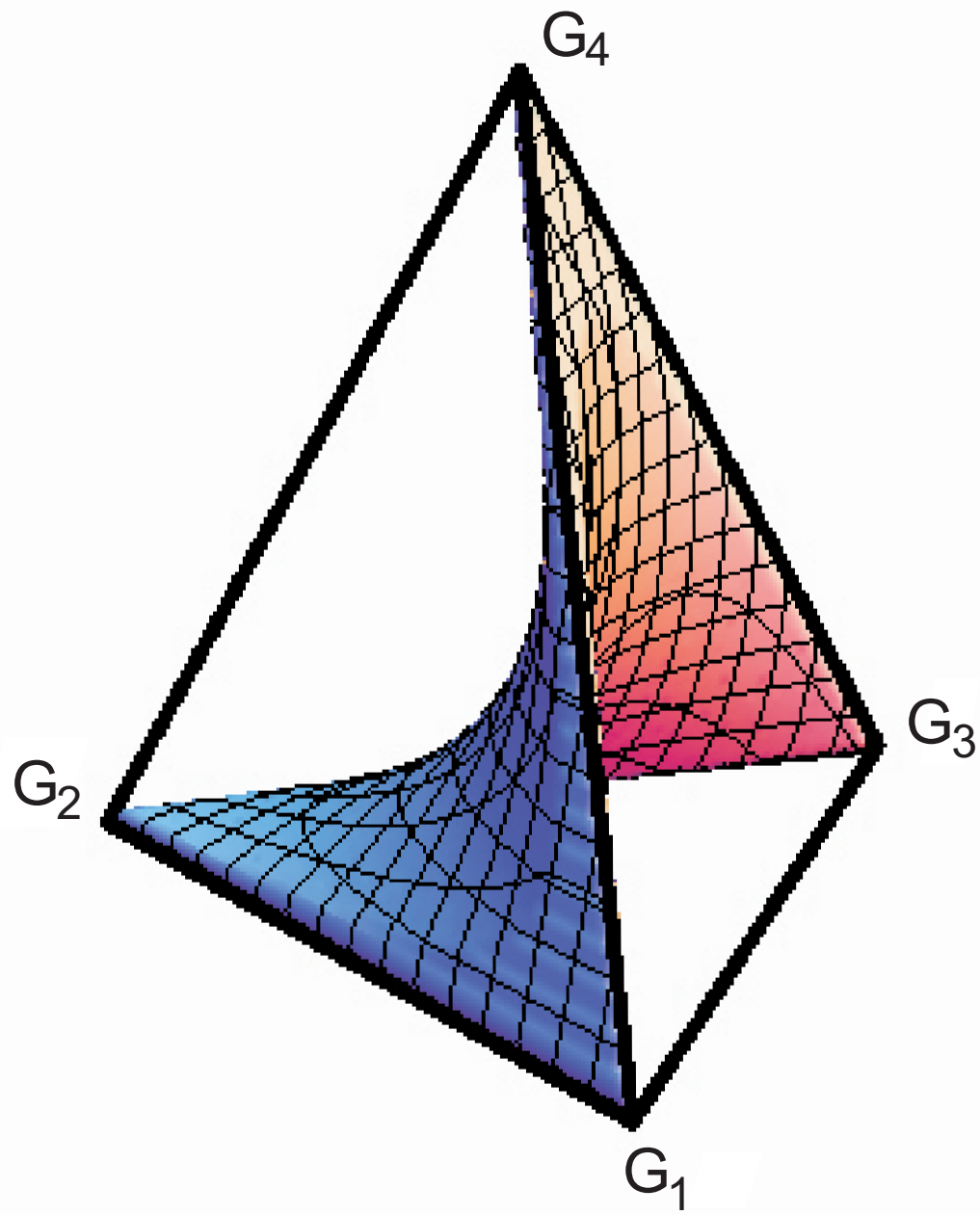
for the own reputation (Haley and Fessler 2005, Bateson et al 2006, Burnham and Hare 2007).
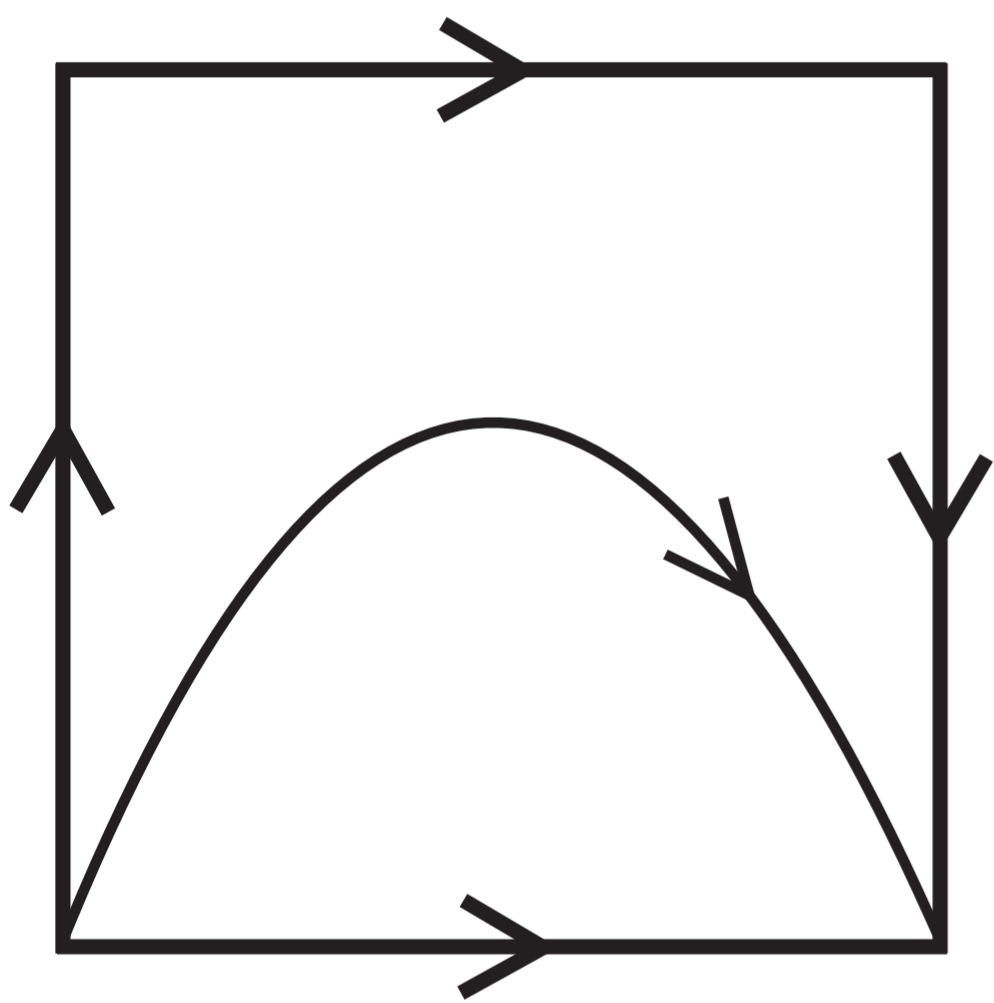
But the emergence of pro-social behavior not only requires information, it also requires a certain amount of selfishness (or 'self-love', to use a kinder but old-fashioned term). Without selfishness, incentives would not work. In the public good games with punishment, for instance, players must not only acquire knowledge about who is a punisher and who not, they must also be prone to defect if they know that they can get away with it. This is a finding well in the spirit of Voltaire's statement that 'it is impossible that a society can emerge and subsist without self-love'.
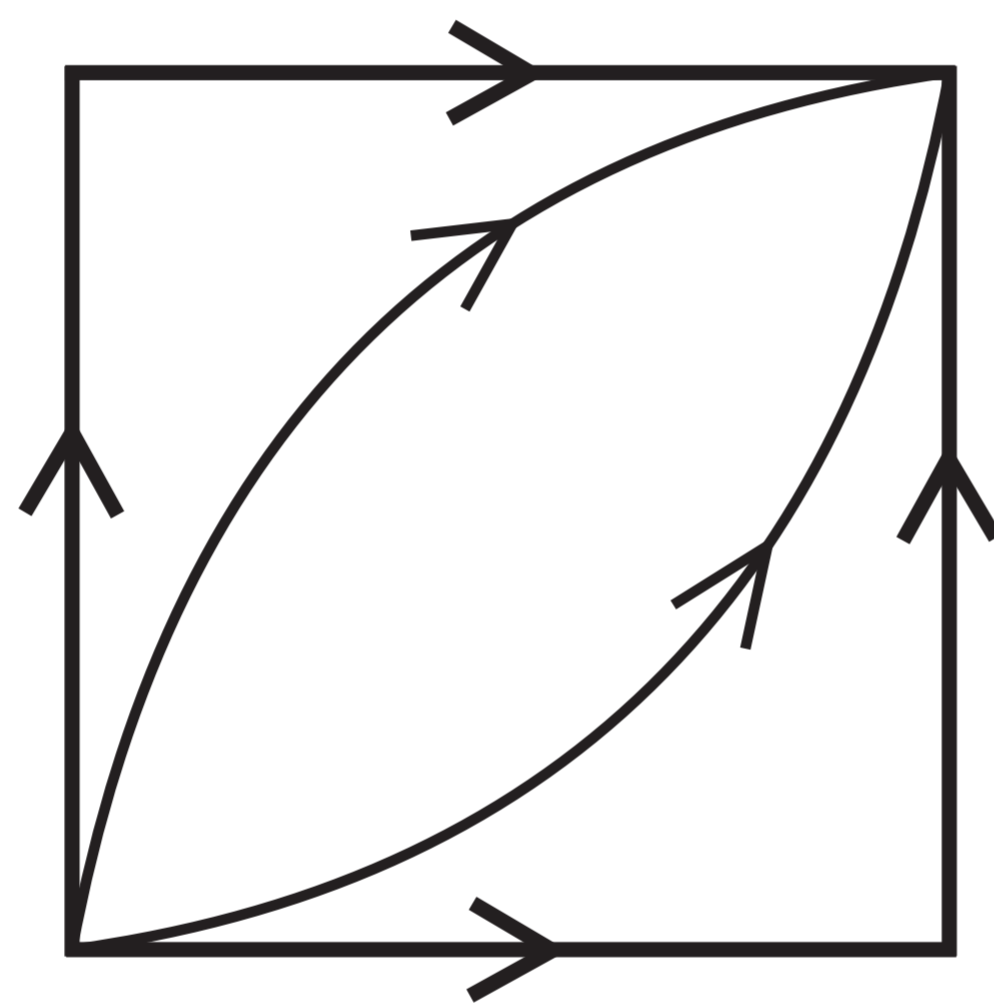
# 15    References

Bateson, M, Nettle, and Roberts, G (2006) Cues of being watched enhance cooperation in a real-world setting, Biol. Lett. 2, 412-414

Berg, J, Dickhaut, J and McCabe K (1995) Trust, reciprocity, and social history, Games and Econ Behav 10, 122-42

Burnham, T and Hare, B (2007), Engineering cooperation: does involuntary neural activation increase public goods contributions?, Human Nature (in press)

Camerer C (2003) Behavioral game theory, Princeton UP

Camerer, C. and Fehr, E. (2006), When does "economic man" dominate social behaviour? Science 311, 47–52

Fehr, E. and Gächter, S. (2002) Altruistic punishment in humans, Nature 415, 137-140

Fehr E and Fischbacher U (2003) The nature of human altruism, Nature 425, 785-791

Güth W, Schmittberger R and Schwarze, B (1982) An experimental analysis of ultimatum bargaining, J Econ Behav Organ 3, 367-388

Haley, K and Fessler D. (2005) Nobodys watching? Subtle cues affect generosity in an anonymous economic game, Evol. Hum. Behav. 26, 245-256

Hardin, G. (1968), The tragedy of the commons, Science 162, 1243–1248

Hauert, C, Haiden, N, and Sigmund, K (2004) The dynamics of public goods, Discrete and Continuous Dynamical Systems B, 4 575-585

Henrich, J. and Boyd, R. (2001), Why people punish defectors, J. Theor. Biol. 208, 7989

Henrich J (2006) Costly punishment across human societies, Science 312, 176-177

Kagel. J H and Roth, A E (eds) (1995) The handbook of experimental economics, Princeton UP, Princeton

Milinski, M, Semmann, D and Krambeck, H J (2002) Reputation helps solve the Tragedy of the Commons, Nature 415, 424-426Nowak, M A and Sigmund, K (2005) Evolution of indirect reciprocity, Nature 437 (2005), 1292-1298

Nowak, MA, Page, K, and Sigmund, K (2000) Fairness versus reason in the Ultimatum Game, Science 289, 1773-1775

Nowak M A (2006) Evolutionary dynamics, Harvard U P, Harvard

Olson M. (1965) The Logic of Collective Action, Harvard University Press

Ostrom, E. and Walker, J. (2003) Trust and Reciprocity: Interdisciplinary Lessons from Experimental Research, Russel Sage Funds

Sigmund, K, Hauert C and Nowak, M A (2001), Reward and punishment, Proc.Nat.Acad.Sci. 98, 10757-10762

Sigmund, K (2007), Punish or Perish? Retaliation and Collaboration among humans, Trends in Ecology and Evolution 22 593-600

Wedekind, C. and Milinski, M. (2000) Cooperation through image scoring in humans, Science 288, 850-852

Xiao, E. and Houser, D. (2005) Emotion expression in human punishment behaviour, Proc. Natl. Acad. Sci. U.S.A. 102, 7398-7401

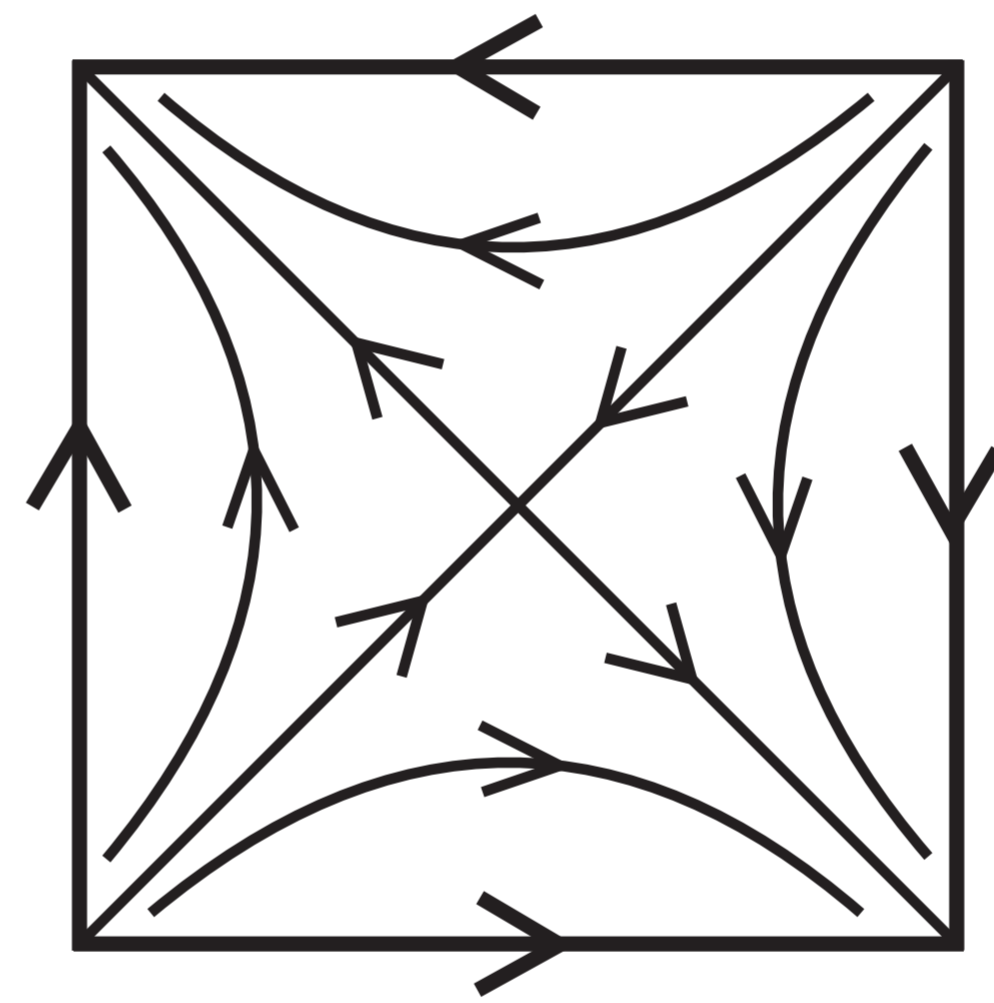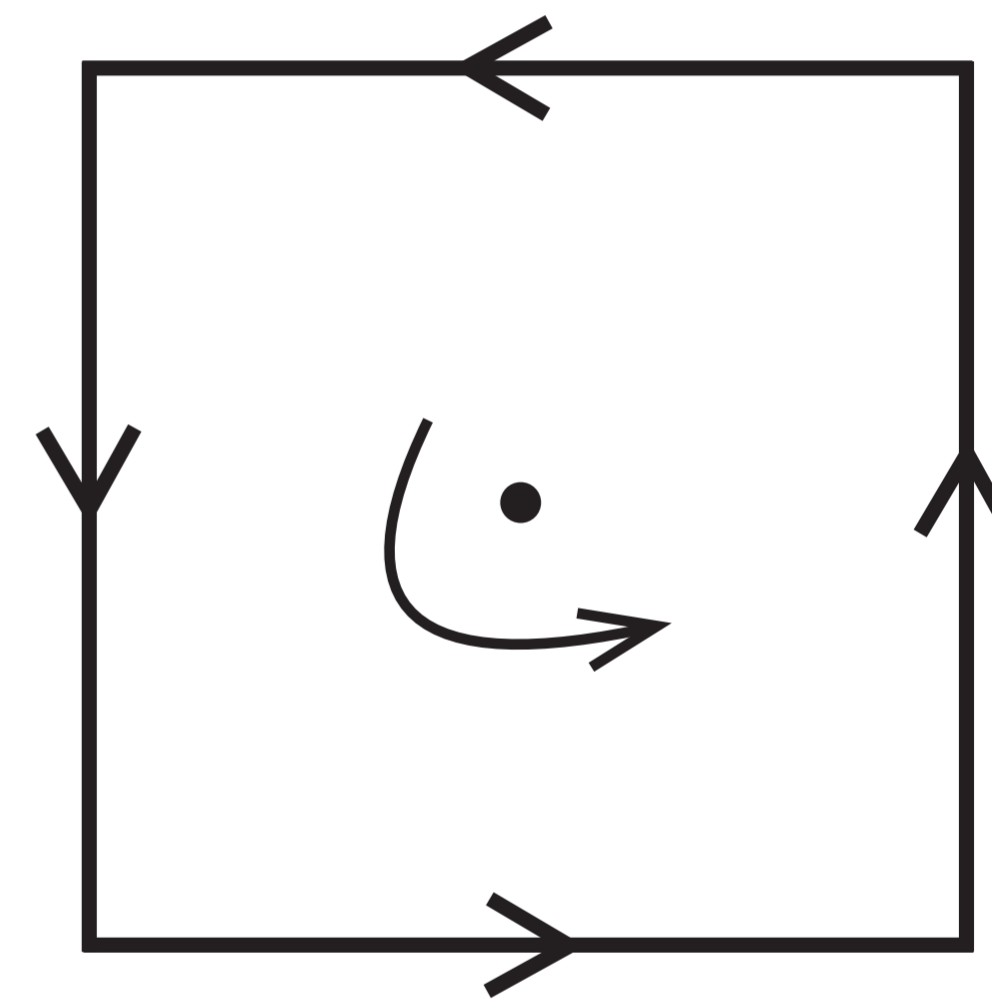a)          b)          c)          d)

mild    $G_4$                       $G_3$   selfish

social    $G_1$                       $G_2$   paradoxical

$G_4$            $G_3$   asocial

pro-social   $G_1$         $G_2$