

Why Punish? Social Reciprocity and the Enforcement of Prosocial Norms

by

Jeffrey Carpenter & Peter Matthews

April, 2002

MIDDLEBURY COLLEGE ECONOMICS DISCUSSION PAPER NO. 02-13



DEPARTMENT OF ECONOMICS
MIDDLEBURY COLLEGE
MIDDLEBURY, VERMONT 05753

<http://www.middlebury.edu/~econ>

Why Punish? Social Reciprocity and the Enforcement of Prosocial Norms*

Jeffrey P. Carpenter
Department of Economics
Middlebury College
Middlebury, VT 05753
(jpc@middlebury.edu)

Peter H. Matthews
Department of Economics
Middlebury College
Middlebury, VT 05753
(pmatthew@middlebury.edu)

April 8, 2002

Abstract: Recently economists have become interested in why people who face social dilemmas in the experimental lab use the seemingly incredible threat of punishment to deter free riding. Three theories have evolved to explain punishment. We survey each theory and use behavioral data from surveys and experiments to show that the theory called *social reciprocity* in which people punish norm violators indiscriminately explains punishment best. We also show that social reciprocity can evolve in a population of free riders and contributors if the initial conditions are favorable.

Keywords: Social Dilemma, Public Good, Punishment, Reciprocity, Norm, Evolutionary Game Theory, Experiment

JEL codes: C91, C92, D64, H41

* We thank Carolyn Craven and Corinna Noelke for comments, Middlebury College for financial assistance, and Carpenter acknowledges the support of the National Science Foundation (SES-CAREER 0092953).

Introduction

Economists have become interested in analyzing, in the experimental lab, something that has been known to field researchers for quite a while – people who face social dilemmas (i.e. situations in which group and individual incentives are at odds) control free riding locally by the use of social, economic, and/or physical sanctions.¹ The existence of schemes by which people monitor each other and punish those who free ride is problematic for standard economic theory. *Why?* First, in non-repeated interactions, any theory assuming that agents simply want to maximize their material gain can not reconcile the cooperative behavior needed to obtain socially efficient outcomes because free riders always do better – free riding when others contribute avoids the costs associated with contributing yet returns the benefits of cooperation and free riding when others free ride prevents one from being taken advantage of. Hence, no selfish person would ever cooperate. The same argument can also be made for not punishing free riders because punishment, in this context, is just a second order social dilemma (see Boyd and Richerson [1992]). Those people who don't punish avoid the costs of doing so, but share any benefits associated with the punishment inflicted by others.

Recently, three theories have emerged to explain why players punish free riders in social dilemma games. In this paper we explain each theory and use data from experiments and surveys to evaluate their behavioral relevance. The common starting point for each theory is the fact that punishing behaviors are grounded in evolutionary theory to provide them with microfoundations. However, the theories differ in the degree to which punishment is believed to be a purposeful act versus a normative response.

Price et al. [forthcoming], hereafter PCT, have developed a theory to explain why people punish free riders. Their theory, which we will call the *fitness differential* theory, differs from standard economic theory because behavior is allowed to evolve based on payoff or fitness differentials. In practice,

¹ Economic *punishment* experiments include Fehr and Gaechter [2000], Bowles et al. [2001], Carpenter [2001], Carpenter et al. [2001], Page and Putterman [2000], and Sefton et al. [2000]. Relevant field research is summarized in Ostrom [1990] and Ostrom et al. [1994]. A specific example of field research is Acheson [1988].

evolutionary models often predict different outcomes from standard economic models because the evolutionary equilibrium selection mechanism is less restrictive. Specifically, non-subgame perfect behavior can be the limit point of an evolutionary process (a la Gale et al. [1995]). The basic premise of the fitness differential theory is the seemingly obvious assertion (as PCT point out) that no pro-social behavior can evolve in a population of free riders unless cooperative types somehow recover or eliminate the benefits of free riding.² Therefore, punishment, according to PCT, evolves to tax away the benefits accruing to free riders.

There are two other theories not based explicitly on removing free riders' fitness advantage that can explain punishment. Both of these theories are based on the notion of reciprocity which requires people to reciprocate the costs imposed by free riders by punishing them. Reciprocal behaviors are similar in phenotype to the behavior posited by PCT (free riders are punished and their fitness is reduced), but the reasons for acting are different. People motivated to reduce fitness differentials are outcome-oriented while, we believe, reciprocators respond to the violation of an obvious norm or rule.

One reciprocity-based theory, originating in Gintis [2000], shows that a behavior called *strong reciprocity*, which is a predisposition to cooperate and punish free riding within well-defined groups, can evolve when group selection is permitted. A second alternative, which we call *social reciprocity* (see Carpenter et al. [2001]) generalizes the notion of strong reciprocity to account for cooperation and punishment in the sort of large, amorphous, groups that constitute neighborhoods, for example. We devote the rest of this paper to assessing how well each theory explains the behavioral data.

Problems with the Fitness Differential Theory in Mutual Monitoring Environments

We define a *mutual monitoring* regime as a system by which people who face social dilemmas in a group with well-defined membership identify and sanction

² It is not actually true that free rider benefits need to be taxed away for cooperation to evolve. For example, assortative interactions allow for the evolution of cooperation by restricting access to the gains from cooperation.

free riders.³ The first question we ask is – *How well does the fitness differential theory predict punishment in mutual monitoring experiments?*

To assess the fitness differential theory we will review the data from an experiment reported on in Carpenter et al. [2001]. In the *mutual monitoring game*, participants are given an endowment of 25 experimental monetary units, EMUs, and given the choice of contributing to a public good or keeping the EMUs for themselves. Each contribution returns 0.5 EMUs for all the members of the four-person group which implies that contributing is socially efficient because each contributed EMU returns 2 EMUs to the group. However, each participant can do better by never contributing because free riders receive 0.5 EMUs per contribution regardless of how much they contribute.

After the contribution stage, players are shown the (anonymous) contributions decisions of the other members of their group and are allowed to punish them. Punishment is costly. Players spend 1 EMU per sanction and each sanction reduces the target's payoff by 2 EMUs.

Instead of declining steadily as usually happens in this sort of choice environment (see Ledyard [1995]), when punishment is allowed, contributions typically start at approximately half the endowment and grow to an average of 70% of the endowment by the end of eight periods.⁴

According to the fitness differential theory, which states that cooperators should punish free riders to reduce their payoff advantage, when looking at the punishment data we should see a strictly downward sloping relationship between contributions to the public good and the punishment assigned. The more one free rides, the larger is that person's payoff differential, and therefore, the more this person should be punished. What we actually see, i.e. the statistically estimated relationship, is given in figure 1.⁵

³ Classic examples in economics are team production and the provision of a public good.

⁴ However, there is typically an end game effect that happens in the last two periods of the game (periods 9 and 10) in which free riding is more frequent, but not punished less.

⁵ Because the data is a panel, the regression behind figure one includes time period random effects. Further, the coefficients on both the first and second order regressors are highly significant ($p < 0.01$). Specifically, we get:

$$\text{Punishment} = 2.12 - 5.22 \text{ Contribute} + 3.23 \text{ Contribute}^2 + \varepsilon$$

(0.26) (0.99) (0.83)

with overall $R^2 = 0.01$, $n=3528$, and Wald chi-squared = 54.91.

Figure 1 here

The horizontal axis of figure 1 measures the fraction of one's endowment contributed and the vertical axis measures the estimated punishment for each contribution level. Because the squared contribution term is highly significant and large ($p < 0.01$, $\beta = 3.23$), the expected punishment function turns up near the end of the contribution range implying that people who contribute at high levels get punished too. This fact is inconsistent with the fitness differential theory – only free riders should be punished.

Given figure 1 is inconsistent with a strict reading of the fitness differential theory, what happens if we consider a less strict version of the PCT theory. For example, what if we say that contributors punish free riders, but free riders also seek revenge and punish contributors. This modification, more or less, still captures the essence of the PCT theory but may be more realistic given other demonstrated tendencies toward spitefulness (e.g. Camerer and Thaler [1995]).⁶

The first problem with this modification is that we need to assume that free riders anticipate being punished by contributors, but a bigger problem is that, on average, free riders don't punish contributors as much as other players do. *Who punishes the contributors?* If we limit our data to the instances in which full contributors (i.e. contribute=1) are punished and regress the total punishment received by full contributors on the contribution level of the person who is doing the punishing, we get the relationship in figure 2.⁷

There is a striking similarity between figures 1 and 2, namely, participants can expect to minimize the amount of punishment they receive by contributing between 70 and 80 percent of the endowment (figure 1), and the people who punish those that contribute everything contribute close to 70 percent (figure 2) themselves. On top of this, the average contribution, pooling

⁶ Also notice that if free riders purposefully seek fitness advantages over contributors then punishing contributors may re-establish the fitness differential eliminated by contributor punishment (if the cost of punishing is less than the harm inflicted).

⁷ Again, we use random effects and get:

$$\text{Punishment} = 1.34 + 1.19 \text{ Punisher's Contribute} - 1.27 \text{ Punisher's Contribute}^2 + \epsilon$$

(0.16) (0.67) (0.71)

with overall $R^2 = 0.04$, $n=76$, and Wald chi-squared = 3.26.

across all periods, is approximately 70% of the endowment and we note that the mean contribution level of the people who punish full contributors is not significantly different from 70% of the endowment ($p=0.95$). These facts put the fitness differential theory in further doubt because even the modified version can't explain the punishment patterns we see. Contrary to the fitness differential theory, contributors are punished and this behavior can not be explained by free riders retaliating against contributors. Instead, players who fully contribute are punished by the players who contribute near the norm or average.

Figure 2 here

If the fitness differential theory doesn't explain the data, how can we make sense of what happens in this game? We think that our mutual monitoring data is more consistent with strong reciprocity if we are more specific about what it means to cooperate. Judging by the data, contributing fully is not necessarily the threshold below which people are determined to be free riders. Instead, as in Fehr and Gaechter [2000] and Bowles et al. [2001], people who deviate from the average are the ones who are punished most. This fact implies that cooperating, in this environment, means conforming to the average contribution level and therefore a modified story in which strong reciprocators will conform to the average and punish people who deviate explains the data rather well.

If our modified strong reciprocity hypothesis is true we would also expect that the likelihood that someone punishes a free rider will depend on how strongly the punisher feels a norm has been violated. To test whether deviating from a norm triggers punishment, we conducted a survey similar to PCT. In this survey, participants answered questions about three vignettes that described a team production scenario in which someone free rides. One relevant scenario participants reacted to is the following:

You and a number of other newly hired people are employed by an auto manufacturer and assigned to work in teams of four. Everyone on the team is paid equally and the pay level is determined entirely by how many cars your work team produces. On the first day of work, you and the other three members of your team divide up the production tasks equally. Over the course of the next month, you and two other members of your group work regularly and hard. However, the fourth member of the team often hides in a storage room and reads

team produces. On the first day of work, you and the other three members of your team divide up the production tasks equally. Each of you works equally hard making cars. However, you notice that a member of the group that occupies the workspace next to yours often hides in a storage room and reads a book instead of working on cars. While your earnings are unaffected by what this member of the other team is doing, the members of his team must work harder and share their income with this person.

We again asked the degree to which the respondent would punish the free rider and the extent to which free riding was considered a violation of a work norm. Surprisingly, 43% of respondents said they would punish the free rider even though, by free riding the person inflicted no harm on the respondent, and the respondent could never benefit from punishing because any increase in effort would benefit a different work team. Further, using this scenario we replicated the regression results discussed in the previous session. The coefficient on the norm measure is again positive and significant ($\beta=0.45$, $p<0.05$) controlling for gender ($\beta=0.74$, n.s.) and the number of economics classes taken ($\beta=-0.28$, n.s.).^{9,10}

The results of our second scenario present a problem for both the fitness differential theory and strong reciprocity. Specifically, why would people who behave according to the fitness differential theory punish free riders in other groups? The same question presents a problem for strong reciprocity. The theory of strong reciprocity hypothesizes that people punish free riders within their group and that groups in which this occurs do better than groups in which punishment does not occur because punishing groups elicit more cooperation. Punishing outside ones group, if any thing, reduces whatever differential benefit punishing groups achieve because punishment is costly and the benefits accrue to a different group.

⁹ The ordered logit procedure in this case yielded:

$$\text{Willingness to Confront} = 0.45 \text{ Free Rider Breaks Norm} + 0.74 \text{ Female} - 0.28 \text{ Econ} + \varepsilon$$

$$(0.20) \qquad \qquad \qquad (0.72) \qquad \qquad \qquad (0.22)$$

with pseudo $R^2 = 0.09$, $n=34$, and chi-squared = 6.97.

¹⁰ We also included a third scenario (and balanced the design) which is identical to the second except it adds a sentence that states that one can expect the free rider to retaliate if confronted. We call this the *high cost* scenario. Again, the relationship between norm violation and punishment is significant ($\beta=0.50$, $p<0.05$).

To account for punishing outside one's group, we define *social reciprocity* as the propensity to cooperate and punish deviations from a widely held norm. So far, all the data we have presented are consistent with social reciprocity. Contributing the average in the mutual monitoring game and punishing deviations is the phenotypic expression of social reciprocity within groups. Outside groups, social reciprocators also punish free riders who break obvious rules like not working.

Social reciprocity and strong reciprocity are related concepts. We view social reciprocity as a generalization of strong reciprocity. One way we might make the theoretical link between the two concepts is to imagine that as society developed, populations grew, congregated in larger numbers, and group membership blurred. In this environment, a simpler heuristic evolved in which people conserve on the cognitive costs of evaluating the group membership of free riders and simply punished all norm violators. That is, it is not hard to imagine a set of parameters such that when one complicates the strong reciprocity model by adding that group membership must be calculated at some cost, but sanctions are relatively costless, people who just punish all violators do as well or better.

Social Reciprocity in the Lab

Because we worried that the hypothetical nature of our survey might bias the result that people will punish "outgroup", we ran another public goods experiment similar to the mutual monitoring game accept, in the second game people could punish outgroup as well as ingroup. We call this game the *social reciprocity game*. The game is identical to the mutual monitoring game accept now, in the punishment stage, participants saw the contributions of, and could punish, *all* the other players in the session and each session was composed of two completely separate groups playing simultaneously.

Figure 3 summarizes play in the social reciprocity game. The vertical axis measures two things simultaneously: the fraction of one's endowment contributed to the public good and the fraction of one's gross earnings spent on punishing other players. The most important thing to notice is that outgroup punishment is positive in all ten periods and responds positively to increased free

ridership. On average, players spend 5% of their earnings punishing free riders outside their group.¹¹

Observe, too, that allowing outgroup punishment increases contributions. Contributions rise not only above the level elicited when no punishment is allowed (Standard VCM), but also above the level resulting when participants can punish within groups only (Mutual Monitor). This makes sense because in worlds populated by social reciprocators free riders are punished more often and more severely which reduces the incentive to free ride. But again, it is important to emphasize that the reason for punishment is not to reduce this differential, or increase contributions, it is a byproduct of the enforcement of a pro-social norm.

Figure 3 here

But Does Socially Reciprocal Behavior Survive Selection?

So far we have provided a third explanation for why people punish in social dilemmas, but we need to prove that such behavior could evolve – that is, *can we construct microfoundations for social reciprocity?* The answer is yes and the microfoundations we provide are a strict test for social reciprocity because we don't use group selection arguments (as in the Gintis [2000] theory of strong reciprocity) nor do we change the environment to favor social reciprocity (i.e. players do not incur costs to identify ingroup members).

Imagine a two person version of the social reciprocity game in which players contribute all or nothing to the public good. The normal form of this game appears in figure 4. As in the experiment, the total endowment of the group is 100 EMUs and contributions to the public good return 0.5 EMUs per EMU contributed. Free riding is clearly the dominant strategy when punishment is not allowed. We also match the experiment by letting agents punish each other. At a cost of 1 EMU players can buy a 2 EMU reduction in another players payoff. For simplicity we make both the contribution and punishment decisions binary; players contribute all 50 EMUs or none and they spend either 10 EMUs to punish each free rider or nothing.

¹¹ Outgroup punishment is significantly positive in each period and directed disproportionately at free riders. See Carpenter et al. [2001] for the details.

	Contribute	Free Ride
Contribute	75, 75	37.5, 87.5
Free Ride	87.5, 37.5	50, 50

Figure 4 – A Mini Social Reciprocity Game

At each moment in continuous time, nature randomly assigns four people from a large population to play the game in two groups of two. All contribution choices are revealed, and then, in a second stage, contributors decide (a) whether to punish free riders and, if so, (b) which free riders to punish – ingroup, outgroup, or both.¹² We consider five agent types:

- (1) *Free Riders*: don't contribute and don't punish.
- (2) *Second Order Free Riders*: contribute but never punish.
- (3) *Strong Reciprocators*: contribute and punish ingroup free riders only.
- (4) *Pure Social Reciprocators*: contribute and punish outgroup free riders only.
- (5) *Social Reciprocators*: contribute and punish free riders in both groups.

We also suppose, for the sake of convenience, that contributors who punish cannot “pick and choose”: a contributor who punishes both in- and outgroup players and is matched with three free riders (one ingroup and two outgroup), for example, is assumed to punish all three.

The mini social reciprocity game has multiple Nash equilibria but a unique subgame perfect equilibrium, in which no one contributes and no one punishes. Suppose, however, that participants are “aspiration-based learners” in the sense of Gale et al. [1995]. In this case, the evolution of the population shares associated with the five pure strategies will follow the normalized replicator dynamic:

$$\dot{p}_i = p_i(\pi_i - \bar{\pi}), \quad i = 1, 2, 3, 4, 5$$

¹² Note this game is an extension of the two person “norms game” described in Sethi [1996], which in turn is based on Axelrod [1984]. See also Binmore and Samuelson [1994] and Gueth and Kliemt [1993].

and one could argue that social reciprocity exists to the extent that strategies 4 and/or 5 survive under this (or similar) selection mechanisms.

Intuition suggests that the expected payoffs of (2) through (5) should be a function of the proportion, p_1 , of free riders alone, and the numbers bear this out:

$$\pi_2 = 75 - 37.5 p_1$$

$$\pi_3 = 75 - 47.5 p_1$$

$$\pi_4 = 75 - 57.5 p_1$$

$$\pi_5 = 75 - 62.5 p_1$$

Are there any circumstances under which social reciprocity survives?

Because free riders will sometimes do much worse than punishers:

$$\pi_1 = 27.5 + 22.5 p_1 + 60 p_2 + 40 p_3 + 20 p_4$$

social reciprocity may evolve from certain initial conditions. Figure 5 illustrates the evolution of behavioral types from the initial state in which 40% of the population are free riders and the remaining strategies each comprise 15% of the population. In this case there are too many free riders for punishment to take hold and despite the second order free riders initially doing well by free riding on the punishers, free riders eventually take over the population.

In a balanced initial population, however, the four contributing types expect $\pi_2=67.5$, $\pi_3=65.5$, $\pi_4=63.5$, and $\pi_5=62.5$, and free riders expect just $\pi_1=56$. The mean payoff for the entire population is $\bar{\pi} = 63$, which means that free riders and social reciprocators will fall short of the population mean, and see their numbers diminish, but the shares of the three other types of contributors will rise. The surprise, perhaps, is that from this initial state, numerical calculation of the replicator dynamic reveals that the free riders and *not* the social reciprocators, will be driven to extinction. The evolution of social reciprocity from balanced initial conditions is shown in figure 6.

Figures 5 & 6 here

A second important fact is illustrated in figure 6. Not only do social reciprocators survive selection in a balanced population, so do strong reciprocators and second order free riders. Perhaps one of the most interesting results of this model is that it predicts a polymorphism of contributing strategies and this polymorphism is, more or less, what we see in the lab. In our experiment, about a third of the participants consistently punish both outside and inside their groups, about half punish ingroup only, and the remaining 20% effectively never punish at all.

Concluding Remarks

Norm-based reciprocal explanations of punishment explain behavioral data better than the Price et al. [forthcoming] theory in which punishment evolves to reduce the differential benefit accrued by free riders. Specifically, we have shown that the Gintis [2000] theory of strong reciprocity can better explain punishment within groups, and a generalization of this theory, which we call, social reciprocity can explain why punishment might occur both within and across groups. Painting in broad strokes, understanding punishment and social reciprocity, in particular, is important because it provides the theoretical foundations to answer the puzzle of how cooperation can evolve and be sustained in large scale populations with fuzzy boundaries, like neighborhoods. Another set of data from neighborhoods in Chicago illustrates the phenomenon we wish to understand. Sampson et al. [1997] show that *collective efficacy*, their term for monitoring and punishing rule breakers locally, is a significant predictor of neighborhood-level outcomes. We claim that such a result can only be understood if we allow for social reciprocators who punish free riders indiscriminately.

Expected Punishment

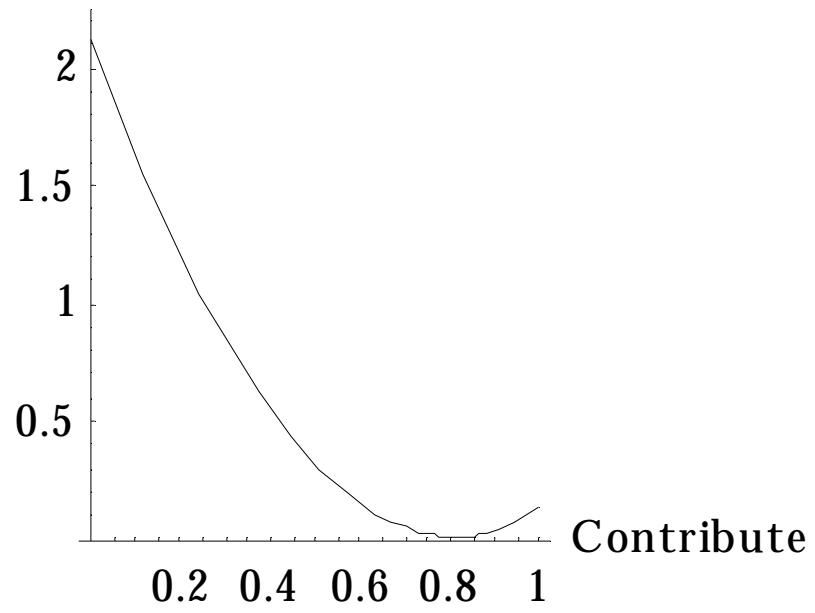


Figure 1 - Who is Punished in the Mutual Monitoring Game?

Assigned Punishment

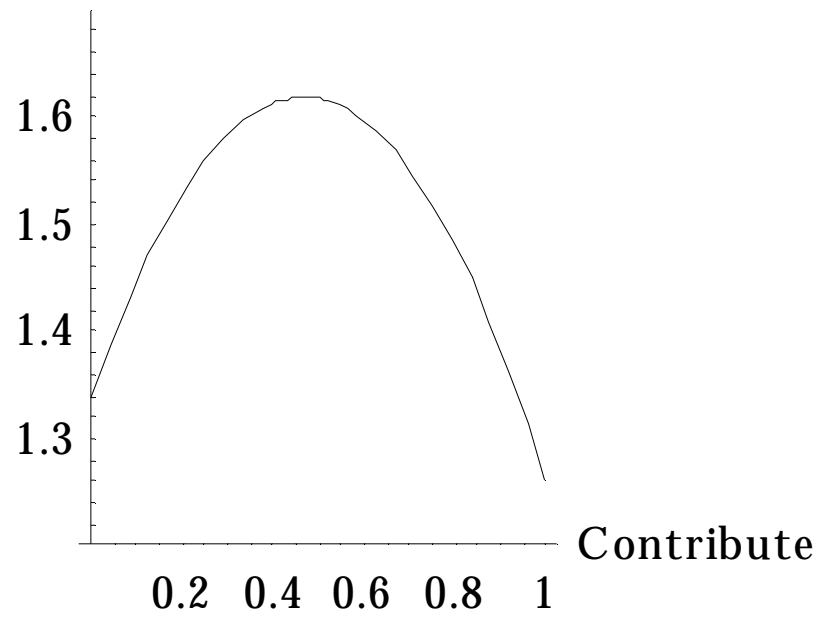


Figure 2 - Who Punishes Contributors?

Contributions and Punishment

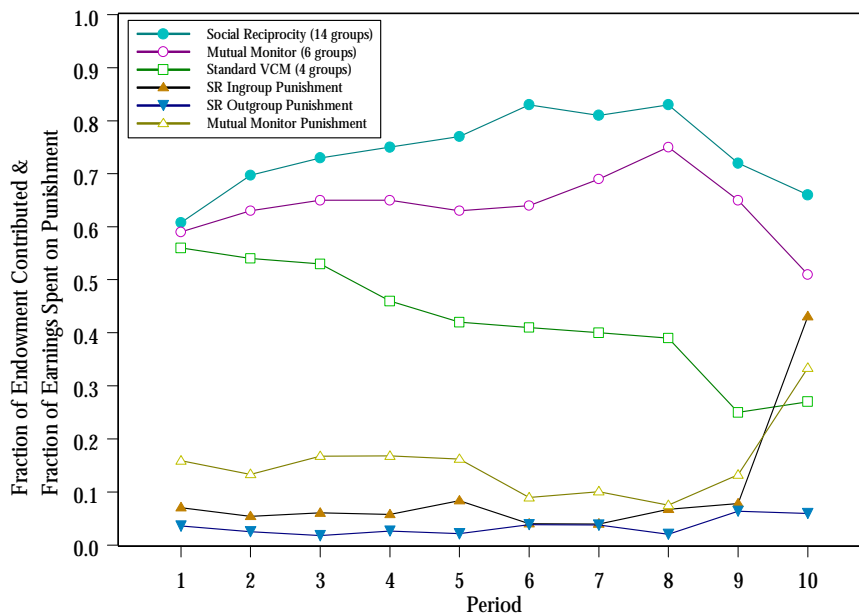


Figure 3 - Behavior in the Social Reciprocity Game

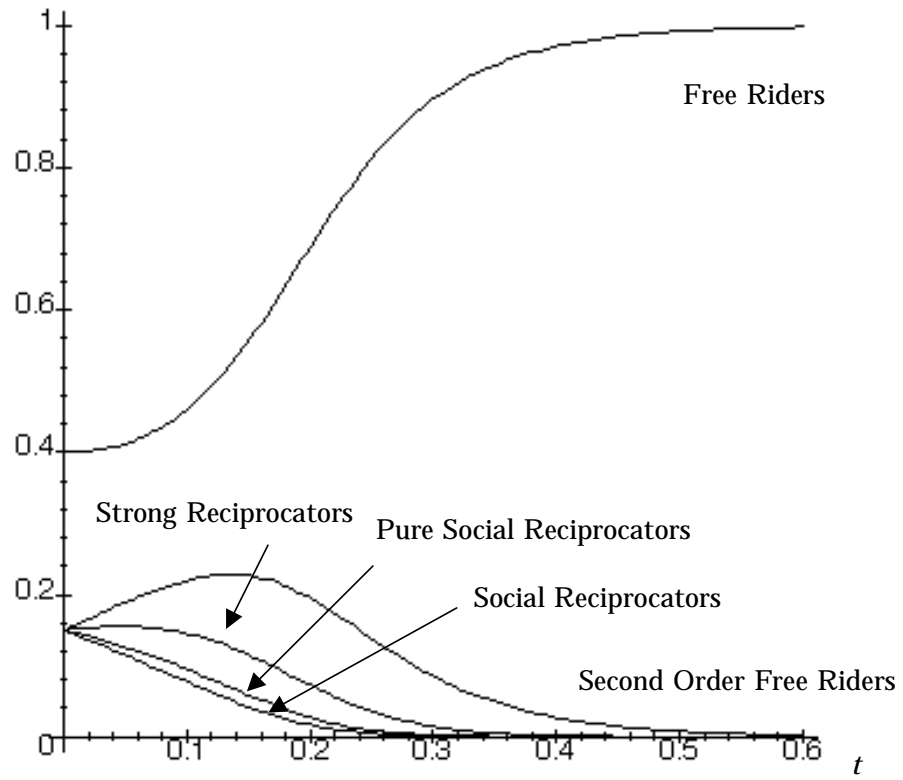


Figure 5 - The Simulated Dynamic with an Unbalanced Initial Population

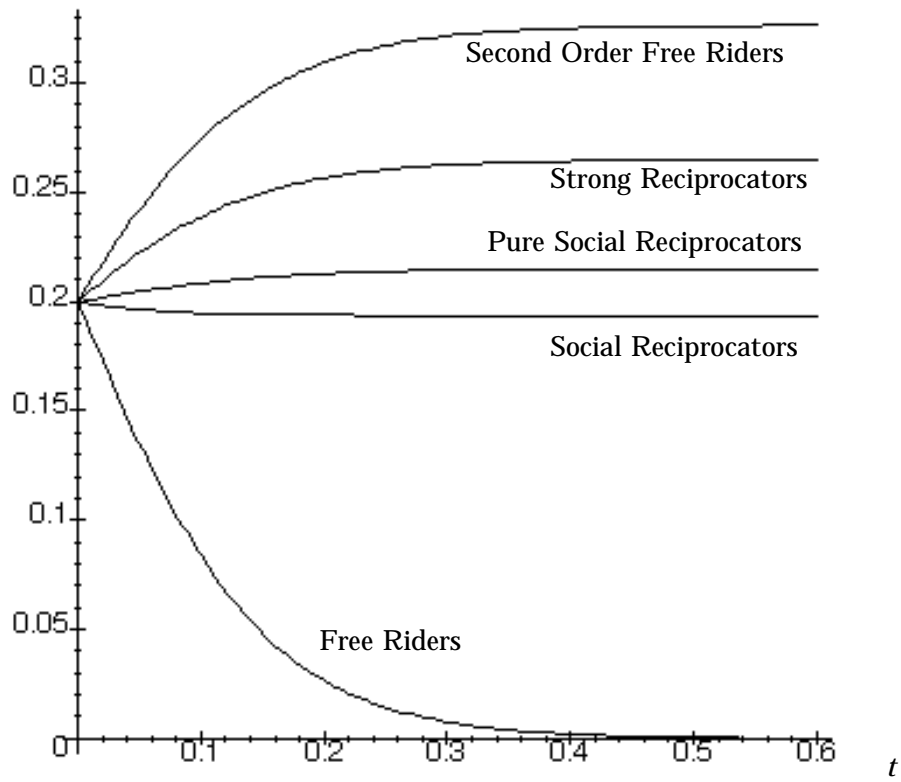


Figure 6 - The Simulated Dynamic from Balanced Initial Conditions

References

- Acheson, J. (1988). *The lobster gangs of Maine*. Hanover, University Press of New England.
- Axelrod, R. (1984). An evolutionary approach to norms. *American Political Science Review* 80: 1095-1111.
- Binmore, K. and L. Samuelson (1994). An economist's perspective on the evolution of norms. *Journal of Institutional and Theoretical Economics* 150(1): 45-63.
- Bowles, S., J. Carpenter and H. Gintis (2001). Mutual monitoring in teams: The effects of residual claimancy and reciprocity. mimeo.
- Boyd, R. and P. Richerson (1992). Punishment allows for the evolution of cooperation (or anything else) in sizable groups. *Ethnology and Sociobiology* 13: 171-195.
- Camerer, C. and R. Thaler (1995). Anomalies: Ultimatums, dictators and manners. *Journal of Economic Perspectives* 9(2, Spring): 209-219.
- Carpenter, J. (2001). Punishing free-riders: How group size affects mutual monitoring and collective action. mimeo.
- Carpenter, J., P. Matthews and O. Ongonga (2001). Social reciprocity. mimeo.
- Fehr, E. and S. Gächter (2000). Cooperation and punishment in public goods experiments. *American Economic Review* 90(4): 980-994.
- Gale, J., K. Binmore and L. Samuelson (1995). Learning to be imperfect: The ultimatum game. *Games and Economic Behavior* 8: 56-90.
- Gintis, H. (2000). Strong reciprocity and human sociality. *Journal of Theoretical Biology* 206: 169-179.
- Gueth, W. and H. Kliemt (1993). Competition or cooperation: On the evolutionary economics of trust, exploitation, and moral attitudes. *Metroeconomica* 45: 155-187.
- Ledyard, J. (1995). Public goods: A survey of experimental research. *The handbook of experimental economics*. J. Kagel and A. Roth Eds. Princeton, Princeton University Press: 111-194.
- Ostrom, E. (1990). *Governing the commons: The evolution of institutions for collective action*. Cambridge, Cambridge University Press.
- Ostrom, E., R. Gardner and J. Walker (1994). *Rules, games and common-pool resources*. Ann Arbor, University of Michigan Press.
- Page, T. and L. Putterman (2000). Cheap talk and punishment in voluntary contribution experiments. mimeo.
- Price, M., L. Cosmides and J. Tooby (forthcoming). Punitive sentiment as an anti-free rider psychological device. *Evolution and Human Behavior*.
- Sampson, R., S. Raudenbush and F. Earls (1997). Neighborhoods and violent crime: A multilevel study of collective efficacy. *Science* 277(August 15): 918-924.
- Sefton, M., R. Shupp and J. Walker (2000). The effect of rewards and sanctions in provision of public goods. mimeo.
- Sethi, R. (1996). Evolutionary stability and social norms. *Journal of Economic Behavior and Organization* 29: 113-140.