

DEPARTMENT OF ECONOMICS

Working Paper

Is altruism bad for cooperation?

by

Sung-Ha Hwang and Samuel Bowles

Working Paper 2008-13



**UNIVERSITY OF MASSACHUSETTS
AMHERST**

Is altruism bad for cooperation?

Sung-Ha Hwang* and Samuel Bowles[§]

31 Aug, 2008

Abstract

Some philosophers and social scientists have stressed the importance for good government of an altruistic citizenry that values the well being of one another. Others have emphasized the need for incentives that induce even the self interested to contribute to the public good. Implicitly most have assumed that these two approaches are complementary or at worst additive. But this need not be the case. Behavioral experiments find that if reciprocity-minded subjects feel hostility towards free riders and enjoy inflicting harm on them, near efficient levels of contributions to a public good may be supported when group members have opportunities to punish low contributors. Cooperation may also be supported if individuals are sufficiently altruistic that they internalize the group benefits that their contributions produce. Using a utility function embodying both reciprocity and altruism we show that unconditional altruism towards other members attenuates the punishment motive and thus may reduce the level of punishment inflicted on defectors, resulting in lower rather than higher levels of contributions. Increases in altruism may also reduce the level of benefits from the public project net of contribution costs and punishment costs. The negative effect of altruism on cooperation and material payoffs is greater the stronger is the reciprocity motive among the members.

JEL codes: D64 (altruism); H41 (public goods)

Keywords: public goods, altruism, spite, reciprocity, punishment, cooperation

Affiliations: * Departments of Economics and Mathematics, University of Massachusetts at Amherst; [§] corresponding author, Santa Fe Institute and Dipartimento di Economia Politica, University of Siena. We thank the Behavioral Science Program of the Santa Fe Institute, the U.S. National Science Foundation, the University of Siena and the European Science Foundation for support of this project and Roland Benabou, Theodore Bergstrom, Simon Gaechter, David Levine, Louis Putterman, Rajiv Sethi, Joel Sobel, and Elisabeth Wood for comments on an earlier draft.

1. Introduction

Both altruism and reciprocity may motivate individuals to contribute to the provision of a public good. Altruism induces the individual to unconditionally value the payoff of other members, while reciprocity implies a valuation of the others' payoffs that is conditional on their contributions (or other indications of their type). Reciprocators may value the payoffs of low contributors negatively and be motivated to reduce the payoffs of defectors at a cost to themselves, when this option is available. The prospect of punishment for low contributions may induce individuals to contribute more than they otherwise would (Fehr and Gaechter (2000), Anderson and Putterman (2006)).

We explore the possibility that these two motives for contribution – a positive valuation of the payoffs of others and a desire to avoid the punishment induced by a negative valuation of one's payoffs by others – may work at cross purposes. Specifically we show that by attenuating the punishment motive, a general increase in the level of unconditional altruism may reduce rather than increase contributions.

Thus, while one often refers to individuals as being 'cooperative' or 'uncooperative', the motives supporting high levels of cooperation are heterogeneous, and they need not work synergistically. For example, experimental evidence indicates that unconditional altruists are significantly less likely to punish low contributors in a public goods game (Carpenter, Bowles, Gintis, and Hwang (Forthcoming)).

In the next section we use the ideas of Levine (1998), Rabin (1993), and Falk and Fischbacher (2006) to explore the joint effects of altruism toward fellow group members and reciprocity-based hostility towards low contributors in a public goods game. In section 3 we study the Nash equilibrium levels of punishment and contribution under varying levels of unconditional altruism of the members of a group. We show that because altruism may diminish the motivation to punish low contributors, the relationship between the level of altruism and contributions is non-monotonic, and that under plausible assumptions there exist a range of levels of altruism over which increases in altruism reduce both equilibrium levels of contribution and the sum of benefits from the public project net of the costs of contributing and the costs of punishing. The latter result is reminiscent of Bernheim and Stark (1988) who showed that increased altruism among two family members in a repeated game setting may be welfare-

reducing. Finally we show that the range for which altruism is bad for cooperation and net benefits is larger the more reciprocal are the group members. In the conclusion we suggest some implications for how social preferences may support cooperation despite the sometimes counterproductive effects of increased altruism and the costly nature of punishment.

2. Altruism, reciprocity and cooperation

Consider a community of individuals indexed by $i = 1, \dots, n$ ($n \geq 3$) who may contribute to a public project by supplying an amount of effort $e_i \in [0, 1]$. The total contributions, $\sum_k e_k$, result in a benefit of $q \sum_k e_k$ which is shared equally among individuals in the community, while each individual experiences the cost of contribution, $1/2 (e_i^2)$. With the notation of $\phi \equiv q/n$, i 's material payoff without the punishment is

$$(1) \quad \pi_i = \phi \sum_k e_k - \frac{1}{2} e_i^2$$

We note that the marginal private benefit of contribution is ϕ and suppose that $1/n < \phi < 1$; $1/n < \phi$ ensures that full contribution, $e_i = 1$, is socially optimal whereas $\phi < 1$ means that in the absence of punishment selfish individuals under-contribute to the public project ($e_i = \phi < 1$).

After contributions have been observed, each individual i can impose a cost on $j \neq i$ with monetary equivalent s_{ij} at cost $c_{ij}(s_{ij})$ to himself. The cost s_{ij} results from public criticism, shunning, ostracism, physical violence, exclusion from desirable side-deals, or another form of harm. Hence $s_i = \sum_{k \neq i} s_{ki}$ is the punishment inflicted upon i by other community members and $c_i = \sum_{k \neq i} c_{ik}(s_{ik})$ is i 's cost of punishing others.

Individual j 's standing as a cooperative member of community, b_j , depends on j 's level of effort and the contribution that j makes to the group, which we assume is public knowledge. Specifically, we assume

$$(2) \quad b_j = 2e_j - 1$$

So $b_j = -1$ if j contributes nothing, and $b_j = 1$ if j contributes fully. This means that $e_j = 1/2$ is the point at which i evaluates j 's cooperative behavior as neither good nor bad. This point could be shifted to any value between 0 and 1, but the added generality is not illuminating.

To model cooperative behavior with social preferences, we say that individual i 's utility depends on his own material payoff π_i , the payoff π_k to other individuals $k \neq i$, the cost of punishing others, and the punishment inflicted on i , according to

$$(3) \quad u_i = \pi_i - c_i - s_i + \frac{1}{n-1} \sum_{k \neq i} (a_i + \lambda_i b_k) (\pi_k - s_{ik})$$

where the parameter a_i , $-1 < a_i < 1$, is i 's level of unconditional altruism if $a_i > 0$ and unconditional spite if $a_i < 0$ and $0 \leq \lambda_i \leq 1$ is the strength of i 's reciprocity motive, valuing j 's payoffs more highly if j conforms to i 's concept of good behavior, and conversely (The function is similar to Levine (1998), but i 's evaluation of k 's type is here based on k 's actions in a particular game, rather than on k 's level of altruism). The valuation of others' payoffs is weighted by the inverse of the number of other members so that changes in group size do not alter the importance of an individual's own payoffs relative to the payoffs of others. The cost to i of punishing j , c_{ij} is increasing in the level of punishment inflicted and it may also increase with i 's level of altruism due to the discomfort that altruists may experience in punishing fellow group members. So we have $c_{ij}(s_{ij}) = 1/2 (a_i + 1)^\kappa (s_{ij})^2$ for $\kappa \geq 0$.

Note (from (3)) that an individual punishing a shirker values the punishment per se rather than the benefits likely to accrue to the punisher if the shirker responds positively to the punishment. Members have an intrinsic motivation to punish the shirker, not simply a desire that the shirker should be punished by someone. This means that punishing is 'warm glow' rather than instrumental towards affecting j 's behavior (Andreoni, 1990; de Quervain et. al. 2004; Casari and Luini, 2008; Anderson and Putterman, 2006). To avoid semantic confusion, note that unconditional altruism and the reciprocity-based spite that motivates punishment of low contributors are both forms of altruism as defined by biologists (assuming that the group benefits associated with the increased contributions induced by punishment outweigh the costs of punishment). Individuals acting according to these motives increase average payoffs in the group

but would enhance their own payoffs were they to (respectively) not contribute or forgo punishing low contributors. We use the term altruism for its unconditional variant.

3. Altruism versus cooperation?

We model a two-stage optimization process in which individual i selects an effort level taking account of the effect of this choice on the punishment inflicted on i by other team members. Because we wish to study the effect of a general increase in the altruism of all group members, we suppose that individuals in the community are homogenous: $\lambda \equiv \lambda_i$ and $a \equiv a_i$ for all i . To find the punishment inflicted on i , we first determine j 's decision concerning the punishment of i depending on i 's contribution level:

$$(4) \quad s_{ji}^*(e_i) = \arg \max_{s_{ji}} u_j(e_j, s_{j1}, \dots, s_{jn}, s_j) \quad \text{for all } j \neq i$$

With $c(s_{ji}) = 1/2(a+1)^\kappa (s_{ji})^2$ member j 's choice of s_{ji}^* in (4) gives the first order condition for an interior solution as follows.

$$(5) \quad c'(s_{ji}^*) = (a+1)^\kappa s_{ji}^* = \frac{1}{n-1} [\lambda(1-2e_i) - a]$$

or the marginal cost of punishing is equal to the marginal benefit of reducing i 's payoffs given j 's assessment of i 's type, net of the subjective costs of inflicting this punishment on i given j 's level of unconditional altruism. When $\lambda = 0$ and $a < 0$, j punishes i , but independent of i 's contribution level. If $\lambda = 0$ and $a \geq 0$, no punishment occurs. If $\lambda > 0$ and

$$(6) \quad e_i \geq e_0 \equiv \frac{1}{2\lambda} (\lambda - a)$$

then member j does not punish. Thus j 's punishment of i is

$$(7) \quad s_{ji}^*(e_i) = \begin{cases} \frac{1}{(a+1)^\kappa (n-1)} [\lambda(1-2e_i) - a] & \text{if } e_i < e_0 \\ 0 & \text{if } e_i \geq e_0 \end{cases}$$

Note that the level of contribution that i must make to avoid punishment by j is declining in j 's level of altruism.

From (7) we can find the total punishment inflicted on individual i , $s_i^*(e_i) = \sum_{j \neq i} s_{ji}^*(e_i)$ which is then non-increasing and differentiable when it is positive. Next individual i decides

the level of effort by taking account of the effect of his effort choice on the level of punishment he will receive. Thus member i will choose

$$(8) \quad e_i(e_{-i}, a) = \arg \max_{e_i} v(e_i) \equiv u_i(e_i, s_{i1}, \dots, s_{in}, s_i^*(e_i))$$

Equation (8) defines member i 's best effort response to other's effort levels, $e_i = e_i(e_{-i}, a)$. To find i 's best response explicitly we proceed as follows. When there is no punishment of i , interior solution of $e_i^{*N}(e_{-i}, a)$ for (8) satisfies the following first order condition (recall $b_j = 2e_j - 1$).

$$(9) \quad e_i^{*N}(e_{-i}, a) = \phi + \frac{1}{n-1} \sum_{l \neq i} (a + \lambda b_l) \phi$$

where $e_{-i} = (e_1, \dots, e_{i-1}, e_{i+1}, \dots, e_n)$

Thus when no punishment is inflicted, i 's optimal choice of e_i equates the marginal cost of contribution (e_i itself) to the direct benefits to i of contributing to the project, ϕ , plus i 's valuation on others' material payoffs. Similarly when i is subject to punishment (hence $e_i < e_0$), i chooses e_i to satisfy the following first order condition :

$$(10) \quad e_i^{*P}(e_{-i}, a) \equiv \phi + \frac{1}{n-1} \sum_{l \neq i} (a + \lambda b_l) \phi - s_i^{*f}(e_i)$$

which requires that i take account of the effect of increased contribution in reducing punishment, as well as the marginal costs and benefits of the project expressed in the no-punishment first order condition (9). Since $s_i^{*f}(e_i) = -2\lambda / (a+1)^k < 0$, we see that $e_i^{*P}(e_{-i}, a) > e_i^{*N}(e_{-i}, a)$; punishment supports a higher contribution level. The amount contributed by i will depend on whether punishment is present or not, and this will depend on the level of unconditional altruism of the members of the group. There exist critical values, \bar{a} and \underline{a} , such that the best response for member i is following.

$$(11) \quad e_i = \begin{cases} e_i^{*P}(e_{-i}, a) & \text{if } a < \underline{a} \\ \frac{1}{2\lambda}(\lambda - a) & \text{if } \underline{a} < a < \bar{a} \\ e_i^{*N}(e_{-i}, a) & \text{if } \bar{a} < a \end{cases}$$

Figure 1 illustrates equation (11).

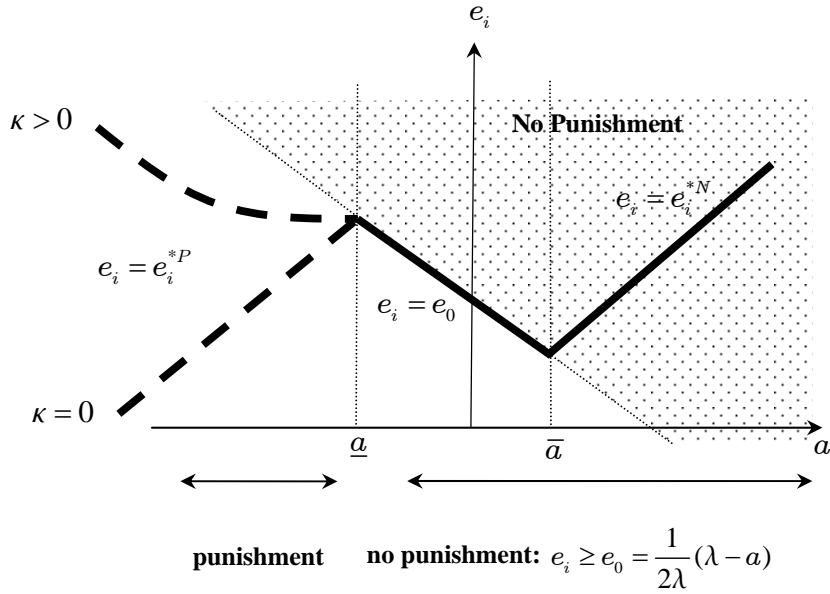


Figure 1. Equilibrium contributions as a function of group member's altruism. To the left of \underline{a} contributions may be rising or falling in a (equation (10)).

When altruism is lower than \underline{a} , i is subject to punishments by others so the effort level is determined by equation (10) and hence may be decreasing or increasing in a . To see this, note that

$$(12) \quad \frac{\partial e_i^{*P}}{\partial a} = \phi - \frac{2\kappa\lambda}{(a+1)^{\kappa+1}}$$

which may be negative if κ and λ are sufficiently large so that the positive effect of altruism (the increased valuation by i of the other members payoffs that are enhanced by i 's contributions) is offset by the negative effect (the increased cost of others punishing i reduces i 's punishment avoidance motives for contributing.) If the cost of punishing does not increase with the individual's altruism ($\kappa = 0$) then $\partial e_i^{*P} / \partial a > 0$, so over this range contributions increase with altruism.

If altruism is greater than \bar{a} , the expected positive effect of altruism occurs because altruism enhances the members' valuation of the external benefits that their contribution allow, while the offsetting effect (the reduced punishment avoidance motive) does not exist because contribution levels are high enough so that punishment does not occur. In the intermediate range

of altruism, equation (6) is binding so an increase in altruism *decreases* the equilibrium effort level since altruism lowers the threshold level of effort required to avoid being punished.

Does the ‘altruism unambiguously bad for cooperation’ range (\underline{a}, \bar{a}) occur for plausible parameter values? Recall that $e_i = \phi$ is the choice of selfish individuals in the absence of punishment and $e_i = 1/2$ is the critical point around which i 's behavior is judged to be good or bad. Thus when the private marginal benefit of contribution, ϕ , is small, so that a selfish individual is motivated to be a bad type (i.e. when $\phi < 1/2$) and members have reciprocal motives (λ is positive), members would punish others and punishment would induce a higher effort level. So we infer that $\phi < 1/2$ and positive λ are necessary conditions for the existence of an interior equilibrium with positive punishment. And if the reciprocity motive is sufficiently strong among community members that the threshold level of effort to avoid punishment, e_0 , reaches 1, an equilibrium with any positive punishment is characterized as full contributions by members. When we exclude cases in which punishment never occurs or in which when it does full contribution is always the result, i.e. when we suppose $\phi < 1/2$ and $1/4 - \phi/2 < \lambda < 0.15$, we obtain the following proposition.

Proposition. We suppose that $1/4 - \phi/2 < \lambda < 0.15$ and $\phi < 1/2$. For $\kappa < \kappa_0$, we can find \underline{a} and \bar{a} such that

$$\frac{de^*}{da} < 0 \text{ for } a \in (\underline{a}, \bar{a})$$

where e^* is a Nash equilibrium. Furthermore, we have

$$\frac{d}{d\lambda}(\bar{a} - \underline{a}) > 0$$

Proof. See appendix. ■

Note that the proposition holds for $\kappa = 0$ (altruism does not affect the cost of punishing). The second part of proposition – that the range over which altruism has a negative effect is increasing in the degree of reciprocity – occurs because the stronger reciprocity motive is, the bigger is the gap between best responses with and without punishment. From $\lambda > 1/4 - \phi/2$ we have $\underline{a} < 0$, so contributions are declining in a not only over the range of positive a but also

over some range of reductions in spite. Note that while increases in altruism for values of a above \bar{a} increase the benefits of the public project net of contribution costs and punishment costs, the reverse is true in the 'altruism unambiguously bad for cooperation' range. Here punishment costs are zero, but increases in altruism reduce contributions to the public good, thus lowering the net benefits. For values of a less than \underline{a} net benefits of the project increase in the level of altruism if contributions also increase. But if contributions are declining in a , then net benefits may either increase or decrease in the level of altruism.

We do not explore the conceptually challenging effect of an increase in altruism on subjective welfare given that the change in altruism is itself a change in preferences (Bergstrom, 2006) analogous to a free resource allowing costless increases in subjective well being. Nor do we address the possibility that were incentive mechanisms other than peer punishment allowed, a general increase in altruism could support more efficient outcomes. If the set of alternative mechanisms is unrestricted this is trivially the case (subsidizing contributions in a complete information setting would achieve this) and there is no non-arbitrary way to expand the set of alternative mechanisms while retaining the underlying problems of public goods provision. Our representation of the motive for punishment - hostility toward those who violate cooperative norms - could be expanded so that the extent of hostility is enhanced by feelings of altruism towards those that the defector has harmed. In this case a general increase in altruism would (as in the current model) make individuals more reluctant to harm defectors, but it would also increase hostility toward defectors, possibly offsetting the first effect. Finally, had we assumed a sophisticated instrumental motive for punishing others, increased altruism could enhance punishment and contributions. The reason is that in this (we think empirically implausible) 'strategic punishing' model, the prospective punisher takes account of the other members' prospective gains resulting from the target's expected positive contribution response to the punishment. For sufficient levels of altruism these gains might outweigh the negative effect of altruism on the non-strategic punishment motive.

Our assumption of a common set of preferences is appropriate for the question we have addressed, but recognizing the heterogeneous nature of preferences would illuminate a further interesting set of issues. In a mixed population of altruists, reciprocators, and self-interested types, for example, reciprocators might punish the altruists as free riders on their civic minded punishment of self-interested defectors. Analysis of the many possible equilibria for this problem

would depend critically on the extent of public and private information and the availability of a common culture or other coordinating mechanisms. We suspect that under plausible assumptions relatively homogeneous sub-populations might outperform mixed populations, and hence might be favored in the process of group formation and the evolution of cultures. But we have not studied this case in detail.

4. Conclusion

Some philosophers and social scientists have stressed the importance for good government of an altruistic citizenry that values the well being of one another. Others have emphasized the need for incentives that induce even the self-interested to contribute to the public good. Implicitly most have assumed that these two approaches are complementary or at worst additive. It is now recognized that this assumption may fail where the presence of monetary or other explicit incentives reduces the salience of altruistic or other public-spirited motives (Benabou and Tirole (2003); Benabou and Tirole (2006); Bowles (2008); Falk and Kosfeld (2006); Sliwka (2007); Bowles and Hwang (2008)). But as we have seen, the assumption need not hold even in the absence of such motivational crowding out.

Our results suggest that for a community wishing to sustain high levels of cooperation, efforts to enhance unconditional altruism may be counter-productive and that enhancing the level of citizen reciprocity may exacerbate the negative effects of altruism. But punishment may also be counter-productive. By definition acts of altruism increase the joint surplus of the community; but punishment is often (as in our model) resource-using. Unless or until levels of contribution sufficient to make punishment rare are achieved, the costs associated with punishment of low contributors may more than offset the gains to cooperation that the punishment allows (Herrmann, Thoni, and Gaechter (2008), Gaechter, Renner, and Sefton (2008)). This is particularly true in a case we have not considered, namely when vendetta-like cycles of punishment and counter punishment are allowed. (Hopfensitz and Reuben (2006)).

Nonetheless, cooperation sustained by a combination of altruism and reciprocity-based punishment may be welfare enhancing. This is true in part because punishment is not only an incentive; it is also a signal. The incentive-based response to punishment may be enhanced by the feelings of shame that punishment by peers triggers (Bowles and Gintis (2006).) In part for this reason disapproval by peers may induce members to contribute even when it is expressed in

non-resource-using ways such as gossip, ridicule or the simple statement that the individual has violated a norm (Maslet, Noussair, Tucker, and Villeval (2003), Barr (2001), Wiessner (2005)).

Appendix

Proof of Proposition

We find $\underline{e} = \underline{e}(a)$ such that $\underline{e} = e_i^{*P}(\underline{e}, a)$.

$$(13) \quad \underline{e}(a) = \frac{\phi}{1-2\lambda\phi}a + \frac{\phi(1-\lambda)}{1-2\lambda\phi} + \frac{2\lambda}{1-2\lambda\phi} \frac{1}{(a+1)^\kappa}$$

We take $\kappa_0 = 3$ and let $\kappa < \kappa_0$ and $g(a) \equiv \underline{e}(a) - 1/(2\lambda)(\lambda - a)$. Then we define \underline{a} satisfying

$$(14) \quad g(\underline{a}) = 0 \quad \text{and} \quad g'(\underline{a}) > 0$$

Such \underline{a} exists since our assumptions ensure that minimum of $g(a)$ achieves at $a < -\lambda$, $g(-\lambda) < 0$, and $g(0) > 0$. Then we verify that \underline{a} satisfies

$$(15) \quad -\lambda < \underline{a}$$

$$(16) \quad e^*(a) < \frac{1}{2\lambda}(\lambda - a) \quad \text{for} \quad -\lambda < a < \underline{a}$$

Similarly we find $\bar{e} = \bar{e}(a)$ such that $\bar{e} = e_i^{*N}(\bar{e}, a)$.

$$(17) \quad \bar{e}(a) = \frac{\phi}{1-2\lambda\phi}a + \frac{\phi(1-\lambda)}{1-2\lambda\phi}$$

We define \bar{a} for which $\bar{e}(\bar{a}) - 1/(2\lambda)(\lambda - \bar{a}) = 0$.

$$(18) \quad \bar{a} = \lambda(1-2\phi)$$

By our choices of \underline{a} and \bar{a} , we have $-1 < -\lambda < \underline{a} < \bar{a} < \lambda < 1$. Now if $a < \underline{a}$ then

$\underline{e} < 1/(2\lambda)(\lambda - a) = e_0 < 1$. Hence when $a < \underline{a}$, $e^* = \underline{e}$ constitutes a Nash equilibrium.

Then for $a > \bar{a} > 0$, $\bar{e} > 1/(2\lambda)(\lambda - a) = e_0$. Thus $e^* = \min\{\bar{e}, 1\}$ is a Nash equilibrium. Finally if $\underline{a} < a < \bar{a}$, then $\bar{e} < e_0 < \underline{e}$ thus $e^* = 1/(2\lambda)(\lambda - a)$ becomes a Nash equilibrium. From this the first part of proposition follows. We summarize this result.

$$(19) \quad e^* = \begin{cases} \underline{e}(a) & \text{if } -\lambda < a < \underline{a} \\ \frac{1}{2\lambda}(\lambda - a) & \text{if } \underline{a} < a < \bar{a} \\ \min\{\bar{e}(a), 1\} & \text{if } \bar{a} < a < \lambda \end{cases}$$

Concerning the second part of proposition, by differentiating $g(\underline{a}) = 0$ with respect to λ we find

$$(20) \quad \frac{d\underline{a}}{d\lambda} = -\frac{-\underline{a}(1-4\lambda\phi) + 2\lambda^2[2(1+\underline{a})^{-\kappa} - \phi + 2\phi^2]}{g'(\underline{a})[2\lambda^2(1-2\lambda\phi)^2]}$$

Since $\underline{a} < 0$ and $\phi < 1/2$ imply $2(1+\underline{a})^{-\kappa} - \phi + 2\phi^2 > 0$, the numerator of (20) is positive. Thus we have $d\underline{a}/d\lambda < 0$. Since $d\bar{a}/d\lambda > 0$, the result follows. ■

References

- Andreoni, James.** 1990. "Impure Altruism and Donations to Public Goods: A Theory of Warm-Glow Giving." *Economic Journal*, 100: 464-77.
- Anderson, Christopher, and Louis Putterman.** 2006. "Do non-strategic sanctions obey the law of demand? The demand for punishment in the voluntary contribution mechanism." *Games and Economic Behavior*, 54: 1-24.
- Barr, Abigail.** 2001. "Social dilemmas, shame-based sanctions, and shamelessness: experimental results from rural Zimbabwe." Centre for the Study of African Economies Working Paper WPS/2001.11: Oxford University.
- Benabou, Roland, and Jean Tirole.** 2003. "Intrinsic and extrinsic motivation." *Review of Economic Studies*, 70 (33): 489-520.
- Benabou, Roland, and Jean Tirole.** 2006. "Incentives and Prosocial Behavior." *American Economic Review*, 96 (5): 1652-78.
- Bergstrom, Theodore.** 2006. "Benefit-Cost in a Benevolent Society." *American Economic Review*, 96 (1): 339-51.
- Bernheim, Douglas, and Oded Stark.** 1998. "Altruism within the Family Reconsidered: Do Nice Guys Finish Last?" *American Economic Review*, 78 (5):1034-1045.
- Bowles, Samuel.** 2008. "Policies designed for self interested citizens may undermine "the moral sentiments:" evidence from experiments." *Science*, 320 (5883):1605-1609.
- Bowles, Samuel, and Herbert Gintis.** 2005. "Pro-social Emotions," in *The Economy as a Complex Evolving System III: Essays in Honor of Kenneth Arrow*, ed. Lawrence E. Blume and Steven N. Durlauf, 337-67, Oxford: Oxford University Press.
- Bowles, Samuel, and Sung-Ha Hwang.** 2008. "Social Preferences and Public Economics: Mechanism design when preferences depend on incentives." *Journal of Public Economics*, Vol 92 (8-9): 1811-20.
- Carpenter, Jeffrey, Samuel Bowles, Herbert Gintis, and Sung-Ha Hwang.** Forthcoming. "Strong Reciprocity and Team Production: Theory and Evidence." *Journal of Economic Behavior and Organization*.
- Casari, Marco and Luigi Luini.** 2008. "Peer Punishment in Teams: Expressive or Instrumental Choice." Unpublished.

- de Quervain, Dominique, Urs Fischbacher, Valeris Treyer, Melanie Schellhammer, Ulrich Schnyder, Alfred Buck, and Ernst Fehr.** 2004. "The Neural Basis of Altruistic Punishment." *Science*, 305: 1254-58.
- Falk, Armin, and Michael Kosfeld.** 2006. "The Hidden Costs of Control." *American Economic Review*, 96 (5): 1611-30.
- Falk, Armin, and Urs Fischbacher.** 2006. "A Theory of Reciprocity." *Games and Economic Behavior*, 52(2): 293-315.
- Fehr, Ernst, and Simon Gaechter.** 2000. "Cooperation and Punishment in Public Goods Experiments." *American Economic Review*, 90 (4): 980-94.
- Gaechter, Simon, Elke Renner, and Martin Sefton.** 2008. "The long-run benefits of punishment." Unpublished.
- Herrmann, Benedikt, Christian Thoni, and Simon Gaechter.** 2008. "Antisocial Punishment Across Societies." *Science*, 319 (7): 1362-67.
- Hopfensitz, Astrid, and Ernesto Reuben.** 2006. "The importance of emotions for the effectiveness of social punishment." *Tinbergen Institute Working Paper 05-0571* (<http://www.tinbergen.nl/discussionpapers/05075.pdf>).
- Levine, David K.** 1998. "Modeling Altruism and Spitefulness in Experiments." *Review of Economic Dynamics*, 1(3): 593-622.
- Maslet, David, Charles Noussair, Steven Tucker, and Marie-Claire Villeval.** 2003. "Monetary and Non-monetary Punishment in the Voluntary Contributions Mechanism." *American Economic Review*, 93 (1): 366-80.
- Rabin, Matthew.** 1993. "Incorporating Fairness into Game Theory and Economics." *American Economic Review*, 83 (5):1281-302.
- Sliwka, Dirk.** 2007. "Trust as a Signal of a Social Norm and the Hidden Costs of Incentive Schemes." *American Economic Review*, 97 (3): 999-1012.
- Wiessner, Polly.** 2005. "Norm enforcement among the Ju/'hoansi bushmen: A case of strong reciprocity?" *Human Nature*, 16 (2): 115-45.