

## The Viability of Vengeance

Daniel Friedman & Nirvikar Singh\*

\* Department of Economics, University of California, Santa Cruz

We offer a theoretical explanation of negative reciprocity or vengeance, the human desire to harm those who have harmed us. Our model shows how negative reciprocity can be sustained by the coevolution<sup>1-3</sup> of genes that determine the capacity for vengeance and group memes (e.g., social norms) that regulate its expression. The model begins with a standard free rider game that captures, simply and directly, a personal cost incurred to reap social gains. The model shows that a taste for vengeance realigns incentives and supports a socially efficient equilibrium, but that by itself the taste for vengeance is not evolutionarily viable. We then show how groups of individuals can use low-power sanctions (or simply status changes) to enforce a particular norm on the proper degree of vengeance. The main result is that actual behavior typically will fall short of the norm, but selection across groups will adjust the norm so that actual behavior maximizes the fitness of group members.

Vengeance is a powerful human motive. We become angry when someone wrongs us, and often try to harm the culprit in return, even at some personal cost. Vengeance deters opportunistic behavior<sup>4</sup> that otherwise might undermine positive reciprocity, direct or indirect<sup>5</sup>, and thereby supports social cooperation. On the other hand, misplaced vengeance sometimes leads tragic feuds and even genocide, as in Rwanda or the Balkans.<sup>6</sup>

The existence of vengeance is empirically obvious (and confirmed in controlled experiments)<sup>5</sup>, but theoretically mysterious, because vengeance is not individually rational: it is weakly dominated by otherwise similar behavior that shirks on the personal cost. Therefore it is a theoretical puzzle how vengeance ever established itself in the repertoire of human motives, and how it sustains itself. Until the puzzle is solved, theory will offer no guidance on how to regulate vengeance to maximize its social value and to minimize its devastation.

( $v=0$ )	C	D
C	1	-1
D	2	0

**Table 1: Fitness with No Vengeance** Entries denote fitness payoffs to a player choosing the row (C=cooperate or D=defect when the partner chooses the column (again C or D). Thus the personal cost (fitness reduction) to choosing C rather than D is 1 and the social gain (own plus partner's payoff) is also 1.

We begin our analysis the simplest possible social dilemma, written as a symmetric 2-player game (Table

1). The game has a unique equilibrium: each player chooses the dominant strategy D and achieves fitness 0, thus missing the potential gains of 1 for each player.

We now add a punishment technology and a punishment motive, parameterized by its incurred cost  $v$ . We hypothesize that a player can inflict harm (fitness loss)  $h$  on the other player at personal fitness cost  $ch$ . The marginal cost  $c$  in  $(0,1)$  is a constant parameter that captures the technological opportunities for punishing others. Also, inflicting harm  $h$  yields the player a utility bonus of  $v \ln h$  (but no fitness bonus) when he is the victim of the sucker payoff (receiving  $-1$ , while the culprit receives  $+2$ ) and no bonus in other circumstances. Thus the motive is not spite<sup>7</sup> but rather is vengeance for damage personally experienced. The motivational parameter  $v$  is subject to evolutionary forces and captures an individual's temperament, e.g., his susceptibility to anger<sup>4</sup>.

The objective function when victim of a sucker payoff now is  $v \ln h - ch - 2$ . Then  $h^* = v/c$  is the utility-maximizing degree of inflicted damage and  $ch^* = v$  the incurred cost. The game now is as in Table 2, and for  $v > c$  we have a coordination game with two locally stable pure Nash equilibria and an unstable mixed Nash equilibrium at  $s^* < 1$  (Figure 1). Thus the threat of vengeance can deter defection and support fully cooperative, socially efficient behavior (C,C) as a Nash equilibrium.

( $v > 0$ )	C	D
C	1	$-1-v$
D	$2 - v/c$	0

**Table 2: Fitness with Vengeance** For  $v > c$ , the strategy D is not dominant. When population fraction  $s$  plays C, the expected fitness of C is  $W(C) = 1s - (1+v)(1-s)$  and the expected fitness of D is  $W(D) = (2 - v/c)s$ . The two expressions are equal at  $s^* = (1+1/v)/(1+1/c)$ . For  $s < s^*$  the expected fitness is higher for D and play converges to the inefficient (fitness 0) all-D equilibrium, as in the basic game. But for  $s > s^*$  the expected fitness is higher for C and play converges to the efficient all-C equilibrium.

**The Viability Problem.** The vengeance motive  $v$  itself is subject to evolutionary forces, perhaps slower than those determining the prevalence  $s$  of cooperation, but real forces nonetheless. The expected fitness of a cooperator is  $W(C|s,v) = 2s - 1 - v(1-s)$ , which is a strictly decreasing function of  $v$  for any fixed  $s < 1$ . Only when there are no culprits ( $s=1$ ) is the expected fitness independent of  $v$ . Thus the fitness of player  $v$  is weakly dominated by that of player  $v'$  whenever  $0 < v' < v$ . Assuming that players occasionally encounter culprits, the vengeance preference parameter  $v$  will be driven towards 0 under any plausible evolutionary dynamics. We have a second-order free rider problem, and it seems that vengeance is not viable.

Standard solutions don't work well for this viability problem. Of course, the problem is attenuated for social creatures that form groups of closely related individuals, such as slime molds (relatedness near  $r=1$ ) or ants ( $r=2/3$ ). But we are interested in humans, whose groups typically consist of individuals who are not that closely related (say  $r=0.05$ ). If each individual's  $v$  were observable, then those with higher  $v$  might encounter D-play less frequently<sup>4</sup> and thus maintain equal or higher fitness. This "greenbeard" solution<sup>8</sup> ignores the evolutionary pressure for lower  $v$  individuals to mimic the visible signs of higher  $v$ . As noted below, building a reputation is a bit complicated in our setting. Less standard solutions are discussed elsewhere<sup>9</sup> and include a continuing stream of mutants<sup>10</sup>, ruling out intermediate values of  $v$ <sup>11-12</sup>, and moralistic strategies of extended negative reciprocity<sup>13</sup>.

**Groups and Memes.** During most of our evolutionary history, humans, like other social primates, presumably lived in small groups of individuals. In constructing our model, we assume that the typical person interacts every day with other members of his or her group and also often interacts with others outside the group, but repeat encounters with any particular outsider are sporadic. In this setting, cooperation within the group can be maintained by many forms of reciprocation (positive or negative, direct or indirect) but cooperation outside the group is problematic; see Methods below for discussion.

All known groups of humans maintain memes that prescribe appropriate behavior towards fellow group members and typically prescribe different appropriate behavior towards individuals outside the group<sup>14</sup>. The analysis below focuses on competition among memes that prescribe behavior towards culprits outside the group and towards group members who deviate from that prescription. These memes determine the group's reputation and therefore the fitness its members receive in encounters outside the group.

The success of the meme, as with any other adaptive unit, is measured by its ability to displace alternatives, i.e., by its fitness. There are many distinct mechanisms by which one meme may displace another, ranging from warfare to fashion, but for the most part these mechanisms align with the most fundamental mechanism, enhanced individual fitness. Without necessarily accepting assertions<sup>15</sup> that misalignments always are minor and temporary, our analysis will assume that a meme prescribing a particular degree of vengeance is fitter than existing alternatives when it brings higher average fitness to group members.

Groups affect individual fitness in several ways. As already noted, they provide gains from internal cooperation and (depending on the reputation they carry) some gains from external cooperation. They also

regulate access to scarce resources such as favorable home sites, stored food and marriage partners. Status or prestige within the group affects access.<sup>14</sup>

**Model Elements.** The model has two parameters describing relevant memes:

- $v^n \geq 0$  is the group's normative vengeance level, the prescribed cost group members are supposed to incur when punishing outgroup culprits; and
- $a \geq 0$  is the group's tolerance of deviations from that norm. A deviation  $x=v^n-v$  incurs a status loss  $x^2/(2a)$  and, of course, corresponding status gains by others.

The relevant genes are summarized in one parameter:

- $v^{max} \geq 0$  is the maximum possible taste for vengeance that any meme could induce, given the individual's capacity for anger and his malleability.

Individual characteristics are described by:

- $v \in [0, v^{max}]$  is the actual vengeance cost an individual prefers; we assume it is learned from personal experience within the group.
- $\bar{v} \in [0, v^{max}]$  is the group average degree of actual vengeance cost incurred in punishing outgroup culprits.

Finally, the impact of group reputation is captured in:

- $f(\bar{v}) = \exp(-\bar{v}/b)$  is the frequency with which an individual encounters culprits. It decreases in the group's reputation, summarized in  $\bar{v}$ , and increases in the positive parameter  $b$  representing the hostility of the environment.

**Result 1: The level of vengeance that maximizes fitness of group members is  $v^o = \max\{0, b-2\}$ .** This level would be advantageous for the group to induce in its members, given the hostility  $b$  of the environment.

The derivation is straightforward. Under present assumptions, sanctions against deviators are zero sum within the group, so the group optimum simply trades off the increased fitness cost of higher  $v$  against the increased benefit, encountering fewer culprits.

Of course, this result describes optimal rather than actual behavior and ignores the second-order free rider problem that no individual captures much of the benefit, which is dispersed throughout the group. To predict actual behavior, we now consider a fixed meme summarized by  $v^n$  and  $a$  and examine how individual preferences will adapt.

**Result 2: Individual adaptation drives  $v$  and  $\bar{v}$  toward the individual optimum  $v^* = v^n - a$ , truncated to  $[0, v^{max}]$ .** For given group characteristics  $\bar{v}$  and  $v^n$ ,

the individual fitness cost is proportional to the sum of the direct cost  $v$  and the loss  $(v^n-v)^2/(2a)$  from deviating from the group norm. The sum is minimized at  $v^* = v^n - a$ . Individual fitness is single peaked at  $v^*$ , so adaptive dynamics push the individual's parameter towards this optimum. The optimum is attained as long as the value

is within the feasible range; otherwise  $v^*$  is truncated below at 0 and above at  $v^{max}$ . Assuming that individual adaptation is faster than either memetic or genetic evolution,  $v^*$  is a good approximation of actual preferences of each individual and an even better approximation of their average.

What is the relation between the group optimum  $v^o$  and the individual optimum  $v^*$ ? The group meme, embodied in the parameters  $a$  and  $v^n$ , is subject to selective pressures in the medium run, and group average fitness is single-peaked. Any group whose memes bring actual behavior (near  $v^*$  by Result 2) closer to true optimum ( $v^o$  by Result 1) has a selective advantage. Hence memetic evolution will drive actual behavior towards the optimum, as long as it is feasible. But if  $v^o$  is infeasible (because the environment is so hostile that  $v^o > v^{max}$ ), then there is a genetic selective advantage to increasing  $v^{max}$ . Hence

**Result 3: Coevolution of memes and genes drives actual behavior  $v^*$  toward the group optimum  $v^o$ ; in equilibrium  $v^n = \max\{0, a+b-2\}$ .**

Thus the actual vengeance level is socially efficient in evolutionary equilibrium. Note that memes  $v^n$  that supports this efficient behavior are not the optimum value  $v^o$  but rather exaggerated versions,  $v^n = v^o + a$ .

On longer time scales there can be shifts in the environment  $b$  and in the punishment technology  $c$ . These shifts affect the encounter function  $f$  and hence the group optimum  $v^o$ . Our main conclusion implies that memes (and, when necessary, genes) will adjust under selective pressure so that individual behavior  $v^*$  will track the new group optimum.

**Discussion.** Vengefulness, or a taste for negative reciprocity, is a crucial part of the human emotional repertoire. We model its importance in sustaining cooperative behavior, but highlight an intrinsic free-rider problem: the fitness benefits of vengeance are dispersed throughout the entire group but the fitness costs are borne personally. Evolutionary forces tend to unravel people's willingness to bear the personal cost of punishing culprits. The countervailing force that sustains vengeance is a group norm together with low-powered (and low-cost) group enforcement of the norm. Such memes coevolve with personal tastes and capacities so as to produce the optimal level of vengeance.

The underlying interaction in our model is the simplest possible social dilemma, but our methods extend easily to more complex interactions. It seems straightforward to redefine a culprit as one who harms any group member, not necessarily oneself. More generally, if the interaction took the form of a common pool resource or public goods game, a player would be considered a culprit to the extent that his contribution falls short by an amount  $e$  from the efficient level (or a normative level). The utility bonus then could take the

form  $ve \ln h$ , in which case the total harm to the culprit would be  $Ve/c$ , where  $V$  is the sum of the other participants' vengeance parameters. Results parallel to 1-3 above seem to follow.

Positive reciprocity could be analyzed directly in a similar fashion: preferences could include a utility gain (but no fitness gain) for rewarding a partner's cooperation, and a social norm could impose a fitness loss on deviators. However, since culprits are rare and cooperators are ubiquitous in successful society, the fitness cost of the rewards is excessive when relying entirely on positive reciprocity to sustain cooperation. Negative reciprocity greatly reduces the burden.

We can speculate how our model applies in different societies. The application to hunter-gatherer bands or villagers is clearest; here parameter  $b$  reflects directly the uncooperative tendencies of people from neighboring groups, and  $c$  reflects the opportunities to identify, track down and inflict harm on them. In more highly structured societies, vengeance is often exacted by delegated specialists. The marginal cost  $c$  of vengeance is lower, but still positive so the model remains valid. Here an important shortcoming is that the model takes as exogenous institutions determining the vengeance technology.

What are the empirical implications of our model? Laboratory experiments can distinguish a taste for negative reciprocity from the egalitarian preferences hypothesized recently<sup>8,16-17</sup>. The model's comparative statics are also clear in principle, and testable with anthropological data: norms of vengeance and vengeful behavior should vary systematically with the hostility of the environment, the technology for harming culprits, and the technology for enforcing group norms. If the model is on the right track, one can hope that dysfunctional vengeful behavior (as in the Balkans) might improve in coming decades as the relevant memes evolve.

## Methods

In this section we explain how results generalizing 1-3 arise from less restrictive specifications than those used in the text. For concreteness we assume that the underlying social dilemma is still as in Table 1.

Even without imposing a group structure, one can write the expected fitness advantage to cooperating  $A(p, u) = W(C) - W(D)$  as a function of the individual's probability estimate  $p$  that his partner will choose C and his expectation  $u$  of her vengeance parameter. One can derive  $p$  and  $u$  from a general specification of noisy observables<sup>9</sup>. Figure 2 shows how the fitter choice, C or D, depends on the sign of  $A$ , and how the choice shifts with the individual's own parameter  $v$ . One can derive<sup>9</sup> a smooth, decreasing encounter function  $f(\bar{v})$  from the assumption that everyone chooses according to Figure 2.

One can also develop a theory of group size and the reliability of group reputation from similar considerations<sup>9</sup>. The text shortcut these matters by implicitly obtaining  $p$  and  $u$  from a convenient but arbitrary encounter function  $f(\bar{v})$  together with a given (high) level  $s$  of within group cooperation.

The crucial feature of the group is that it imposes an expected fitness loss  $\rho(x)$  when an individual deviates  $x = v^n - v$  from the group norm  $v^n$ . The loss function  $\rho$  is assumed smooth, convex (with slope  $>1$  for  $x$  sufficiently large) and minimized at 0, but it need not be quadratic or even symmetric.

**Result 2 Derivation.** With probability  $f(\bar{v})$  a  $v$ -cooperator encounters a defector and receives fitness loss  $(1+v+\rho(v^n - v))$ , the sucker payoff plus the cost of wreaking vengeance plus the group status loss from departing from the norm. The individual receives a fitness gain of 1 in an encounter with a cooperator. Thus the individual's expected fitness is  $W(v | \bar{v}, v^n) = 1(1-f(\bar{v})) - (1+v+\rho(v^n - v))f(\bar{v}) + R = 1 - f(\bar{v})(2+v + \rho(v^n - v)) + R$ , where  $R$  is the base-level fitness including the (positive) effect on one's status from other group members' deviations from the norm  $v^n$ . The fitness function does not account for the possibility that the individual may sometimes play D, but this omission is harmless in terms of analyzing the adjustment of  $v$ .

Since the current choice of  $v$  has negligible effect on  $\bar{v}$ , the derivative of  $W$  is, up to the positive multiplicative constant  $f(\bar{v})$ , simply  $\rho'(v^n - v) - 1$ . The assumptions on  $\rho$  ensure that this derivative is decreasing in  $v$  and equal to zero at a unique  $v^* > v^n$ . Hence the single-peaked property holds and the argument in the text shows that individual adaptation pushes actual  $v$  towards  $v^*$ . With the quadratic specification for  $\rho$  we have the first order condition  $0 = \rho'(v^n - v) - 1 = (v^n - v)/a - 1$ , so in this case  $v^* = v^n - a$ , truncated to  $[0, v^{max}]$ , as claimed.

**Result 1 Derivation.** The text assumes that the group imposes  $\rho$  purely through status changes. In this case,  $R$  cancels the mean contribution of  $\rho$ , so the group average fitness is  $W^g(\bar{v}) = 1(1 - f(\bar{v})) - (1 + \bar{v})f(\bar{v}) = 1 - f(\bar{v})(2 + \bar{v})$ . The group optimum  $v^o$  maximizes this expression on  $(c, v^{max}]$ . Using  $f(v) = \exp(-v/b)$ , the first order condition reduces to  $2 + v = -f/f' = b$ , so  $v^o$  is  $b - 2$ , truncated to  $(c, v^{max}]$  as claimed.

More generally, imposing a loss on a deviator may reduce the group's overall fitness by some fraction  $t \in [0, 1]$ ; e.g., some potential gains to cooperation may not be realized. Group average fitness becomes  $W^g(\bar{v}) = 1 - f(\bar{v})(2 + \bar{v} + t\rho(v^n - \bar{v}))$ . With  $f$  and  $\rho$  as specified in the main text, direct computation yields  $v^o = a(-t/2) + b(1 - t) - 2$ . If  $t = 0$ , we have the case just analyzed, where  $v^o$  depends only on the environmental hostility parameter  $b$ .

Higher  $t$  decreases the optimal level of vengeance and introduces positive dependence of  $v^o$  on the tolerance parameter  $a$  and the environmental hostility parameter  $b$ .

**Result 3 Remarks.** The text uses a direct argument from earlier results to establish the striking result that memes solve the second-order free rider problem and align actual behavior with optimal behavior. Here we highlight two underlying assumptions. The first is that there is a hierarchy of time scales, so that individual  $v$  (hence also  $\bar{v}$ ) adapts most rapidly, then the memes  $v^n$  and/or  $a$ , then the genes  $v^{max}$ , and finally the environmental and technological parameters  $b$  and  $c$ . It does seem that  $v$  adapts rapidly to social memes; e.g., according to stories in the media, kids raised in Belfast and Lebanon brought to the US have no problem adapting with a few months to the US norm and then adapting back when they return. It is natural to think of memes as evolving faster than genes, but that is not necessary for our result. Examples of coevolution on overlapping time scales include sickle-cell anemia in yam-growing areas, and lactose tolerance in herding communities<sup>2</sup>. The optimality result clearly fails when the environment or punishment technology changes faster than the memes (or, when  $v^{max}$  is too low, the genes).

The other key assumption is that there is no fitness conflict between memes and genes. Generally speaking, memes transmitted vertically (e.g., from parents to children) share a common fate with genes and hence their fitnesses tend to coincide. However, horizontal meme transmission (e.g., from one teenager to another) need not respect genetic fitness. Fortunately for our argument, available anthropological evidence points to vertical transmission of  $v^n$  and  $a$ .<sup>14</sup>

September 13, 1999

### References

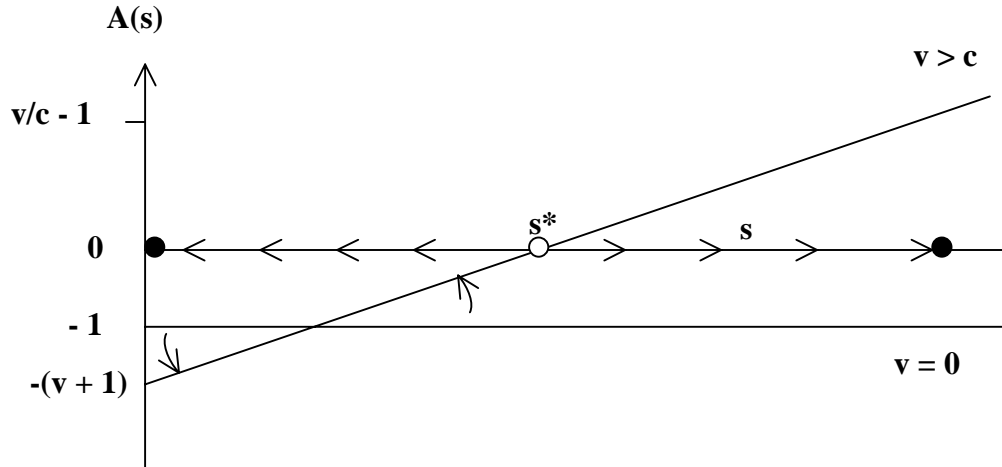
1. Dawkins, Richard (1976), *The Selfish Gene*, NY: Oxford University Press.
2. Durham, William H. (1991), *Coevolution : genes, culture, and human diversity*, Stanford, Calif.: Stanford University Press.
3. Boyd, Robert and Peter J. Richerson (1985), *Culture and the Evolutionary Process*, University of Chicago Press.
4. Frank, Robert (1988), *Passions within Reason: The Strategic Role of the Emotions*, NY: WW Norton.
5. Fehr, Ernst and Simon Gächter (1998), "Cooperation and Punishment," University of Zurich manuscript, September
6. ..New Yorker article of Sept 6, 1999 or other discussion of vengeance in the balkans...
7. David K. Levine (1998) "Modeling Altruism and Spitefulness in Experiments," *Review of Economic Dynamics* 1, 593-622.
8. Boyd & Richerson greenbeard cite..
9. Friedman, Daniel and Nirvikar Singh (1999), "On the Viability of Vengeance", UC Santa Cruz Working Paper, <http://...>
10. Sethi, Rajiv and Eswaran Somanathan (1996), The Evolution of Social Norms in Common Property Resource Use," *American Economic Review*, 86:4, 766-788.
11. Huck, Steffen and Jorg Oechssler (1999), "The Indirect Evolutionary Approach to Explaining Fair Allocations", *Games and Economic Behavior*
12. Robert Axelrod, (1986), "An Evolutionary Approach to Norms," *American Political Science Review*, 80, 1095--1111.
13. Boyd, Robert, and Peter J. Richerson, "Punishment Allows the Evolution of Cooperation (or Anything Else) in Sizable Groups", *Ethology and Sociobiology*, 13, 171-195.
14. Sober, Elliott and David Sloan Wilson (1998). *Onto others: The evolution and psychology of unselfish behavior*, Harvard University Press.
15. Wilson, Edwin O...Sociobiology...
16. Bolton and Ockenfels (1998). "ERC: A Theory of Equity, Reciprocity and Fairness," Penn State University manuscript.
17. Fehr, Ernst and Klaus Schmidt (1999) "A Theory of Fairness, Competition and Cooperation," QJE.

Correspondence and requests for materials should be addressed to D.F. (email: dan@cats.ucsc.edu).

**Acknowledgements.** The first author is grateful to CES Munich University for hospitality while writing the first fragments in May 1997. We have benefited from the comments of audiences at Indiana, Purdue, UCLA, and UCSC.

**Figure 1: The Advantage of Cooperating.**

The fitness advantage  $A(s) = W(C) - W(D) = (v/c - 1)s - (v+1)(1-s)$  is graphed as a function of the population fraction  $s$  playing C for two values of the vengeance parameter  $v$ . The graph of  $A$  rotates counterclockwise as  $v$  increases. The solid dot at  $s=1$  represents a socially efficient, fully cooperative Nash equilibrium.



**Figure 2: The Decision Rule.**

The appropriate choice of C or D is given by the sign of the advantage function  $A(p,u)$ , where  $p$  is the probability that the partner will choose C and  $u$  is an unbiased estimate of her vengeance parameter. The  $A=0$  locus shifts up with increases in the the decision maker's direct ( $v$ ) or full ( $\alpha = v + \rho(v^p - v)$ ) vengeance cost.

