

Aversion to Norm-Breaking: A Model¹

Raúl López-Pérez²

February 2007

¹I am indebted to Jordi Brandts, Marco Casari, Gary Charness, Danilo Coelho, Simon Gächter, Erin Krupka, Clara Ponsatí, Debraj Ray, Andrew Schotter, Angel Solano, Christian Traxler and seminar participants at UAB, NYU, the March 2004 Public Choice meeting and the November 2003 Urrutia Elejalde Foundation and UNED Winter Workshop on Economics and Philosophy for very helpful comments. Part of this research was conducted while visiting NYU and I would like to thank the Economics Department -specially to Andrew Schotter- for its great hospitality. I also gratefully acknowledge financial support from the Spanish Ministry of Education and Science.

²Department of Economic Analysis, Universidad Autonoma de Madrid, 28049 Madrid, Spain.
E-mail: raul.lopez@uam.es

Abstract

In experimental games, we observe the following phenomena: (1) Many subjects cooperate contrary to their material interest, (2) they cooperate in a reciprocal manner, (3) subjects often punish those others who behave unkindly, and (4) previous history usually influences subjects' choices. We propose a simple game-theoretical model to account for these and other experimental phenomena, and compare it with other models of social preferences and reciprocity.

Keywords: Emotions, Fairness, Path-Dependency, Reciprocity, Social Norms. JEL classification numbers: C70, C72, D63, D64, D74, Z13.

1 Introduction

Experimental Economics offers abundant evidence -see Fehr and Schmidt (2006) and Camerer (2003) for excellent surveys- that contradicts the joint hypothesis that *all* agents are rational and motivated *only* by their own material interest. In a Dictator ‘game’ experiment, for instance, one subject is provisionally endowed with some money and must decide how much of that money to transfer to another (anonymous) participant; the ‘game’ finishing then. Clearly, a rational and selfish chooser would not transfer anything. Contrary to that prediction, a significant proportion of the participants give something, many times as much as half of the stake.

The Ultimatum game, another well-known experimental game, provides additional evidence in this line. This game has the same structure as the dictator game except that the second player (the ‘responder’) has now a say and can accept or reject the first mover’s proposal of sharing. The proposal is implemented if the responder accepts it, whereas both players get zero money if it is rejected. Obviously, the rejection of a strictly positive offer goes against material interest. However, actual responders usually reject offers of less than one quarter of the stake and even more.

These results are proof that (some) people cooperate (or behave generously towards others) *and* punish contrary to their material interest. Why do cooperation and punishment occur? The dictator’s decision problem and the responder’s one are so simple that an argument based on rationality failures seems rather convoluted. On the contrary, introspection suggests that motivational forces different from material interest play here a crucial role. This paper investigates formally such motivations in order to offer a rational choice explanation of subjects’ behavior in this and many other experiments.

We believe that the aforementioned phenomena can be explained by resorting to social norms and emotions. Human societies are endowed with social norms that people internalize through the education process. In this way, people acquire certain emotional responses to others’ and one’s behavior. On the one hand, self-conscious emotions like embarrassment, guilt or shame trigger when *oneself* deviates from an internalized norm. On the other hand, people feel aggressive emotions like anger when *another player* violates a norm that one has hitherto respected. These *painful* emotions shape human preferences because, other things constant, one prefers not to suffer them. In turn, norm internalization affects human behavior in two distinctive ways. First, people adjust their

choices to prevent the activation of the above mentioned negative emotions. Second, specific behavioral impulses appear associated with such sensations once they get triggered (Frijda, 1986) -the action tendency of anger, for instance, is to punish the deviator.

As a result, emotions crucially shape norm compliance and punishment in human societies: People respect norms to avoid bad feelings (*internal punishment*), external sanctions (*external punishment*), or both, whereas, in addition, human punishment is often driven by aggressive emotions. Some recent neurological evidence is consistent with the idea that emotions induce punishing behavior. Sanfey et al. (2003), for instance, scanned responders' brains in the ultimatum game and found that rejection of low offers was correlated with heightened activity in the anterior insula (a brain area thought to be related to emotions like anger). In addition, self-reported questionnaires of emotional states suggest that emotions like shame, guilt or anger play an important role in decision making (Bosman and van Winden, 2002; Fehr and Fischbacher, 2004).

From our point of view, the generosity and the Pareto-damaging behavior that (some) subjects exhibit in the Dictator and Ultimatum games, respectively, may be explained by the theory sketched in the two previous paragraphs. The argument is simple: (Some) subjects have internalized a specific norm of fairness or distributive justice which they take to the lab, and then their emotions make them act according to the norm and punish transgressors. Intuitively, these principled subjects also affect the behavior of the remaining, self-interested, agents which may find profitable to respect the norm if they risk being sanctioned otherwise. This paper provides a formalization of this theory.

We analyze the working of the model in different strategic settings like the Dictator game, Cournot game, Ultimatum game, Trust game, Centipede Game, Competition games, and Best-Shot game, and show that it is more consistent with experimental evidence than the standard *homo economicus* model. Although this may be debatable, we also believe that our model has several advantages over other models of social preferences and reciprocity, among which we may cite some. First, our model is very general: One may use it not only to understand why people respect fairness norms, but also why they follow dressing norms, codes of etiquette, or honesty norms, which other models have difficulties to explain. Second, and contrary to *some* other models, it is a model of *path-dependent preferences*: Players care about their material payoff, but also about *how* they get it. This appears to be largely consistent with experimental evidence -again, consult,

Fehr and Schmidt (2006) or Camerer (2003). Third, and consistent with the extensive evidence provided by Charness and Rabin (2002), the model predicts two apparently contradictory phenomena: (i) Many subjects have *both* social efficiency and equality concerns, and (ii) many subjects engage in Pareto-damaging behavior to punish deviators.

Last, but not least, ours is a model of reciprocity in the sense that principled players respect the norm in equilibrium only if they expect sufficiently others to respect it as well, and they hurt transgressors. Further, it is relatively parsimonious because it assumes that players do just care about their co-players' behavior, *not* about their motives for such behavior -i.e., their intentions. This contrasts with other models of reciprocity, which assume that players care about their co-players' beliefs or types and thus have to resort to the Psychological Game Theory of Geanakoplos et al. (1989) or to signalling arguments that complicate much the analysis.

In order to avoid misunderstandings, let us be very clear on one point: We do not have any doubt that intentions are important to explain decision making. Therefore, our model is unrealistic in this regard. However, we believe that this is compensated by the fact that the model is relatively simple and (yet) empirically relevant. This latter feature is remarkable: Our model can explain the kind of experimental phenomena that have been usually associated with intentions, that is, (1) the influence of previously non-chosen alternatives on behavior, and (2) the fact that players discriminate between passive and active co-players *in certain circumstances* and in a way that contradicts the assumption that players only care about the distribution of material payoffs. For this reason, we believe that our paper could contribute new insights to the debate on when intentions affect behavior.

The remainder of the paper is organized as follows. We devote the next section to survey some of the literature on social preferences. To distinguish the effects of prosocial and aggressive emotions on behavior, section 3 first describes prosocial preferences and section 4 applies this model to different experimental games, comparing theoretical predictions with experimental data. Section 5 adds aggressive emotions to the model and studies other experimental games. Throughout sections 4 and 5, we point out the differences between our approach and that of other models. Finally, section 6 proposes some possible extensions and concludes.

2 Other Models of Social Preferences and Reciprocity

However the pervasiveness of the *homo economicus* hypothesis, the idea that people have social emotions has an old history in Economics. Edgeworth (1881) proposed a simple model of altruism in which an individual's utility is a weighted sum of her and others' *material* payoffs. This linear formulation is rich enough to express many ideas. To describe it in game theoretical terms, assume for simplicity that player i 's material payoff at terminal node z coincides with her money earnings $x_i(z)$. Player i 's utility at z is then

$$U_i(z) = x_i(z) + \sum_{k \neq i} \alpha_{ik}(z) \cdot x_k(z), \quad (1)$$

where $\alpha_{ik}(z) \in [-1, 1]$ for any i, k , and z . Of course, the *homo economicus* hypothesis entails $\alpha_{ik}(z) = 0$ for any i, k and z . On the opposite, it is said that player i is altruistic towards player k at z if $\alpha_{ik}(z) > 0$, and spiteful towards k if $\alpha_{ik}(z) < 0$.

In the simplest formulation within this linear framework, $\alpha_{ik}(z)$ is the same constant number for any k and z . This means that the sign and intensity of our sentiments or emotions towards the others do not depend on their acts, qualities, and beliefs, or on the actual distribution of material payoffs.

In an important and pioneering paper, Rabin (1993) put into question the previous formulation, providing an alternative model. Rabin pointed out that the sign of our sentiments is conditional: "[...] *the same people who are altruistic to other altruistic people are also motivated to hurt those who hurt them.*"¹ Moreover, he posited that the sign of our sentiments depends on our beliefs about the others' intentions.

Roughly speaking, player B's intentions are her *expectations* about the terminal distribution of material payoffs to be reached in the game. Take then any two-player game in normal form and suppose that A *believes* that B's intentions are (x_A^*, x_B^*) . B's *intentions* are kind (unkind) to A if x_A^* is larger (smaller) than the equitable payoff -i.e., the average of the maximum and minimum A's payments within the set of Pareto efficient allocations that, according to A, B believes to be reachable. In a somewhat analogous way, B is kind (unkind) to A if she expects A to get a higher (lower) payoff than what B believes to be A's equitable payoff. Now, a player's utility is the sum of her expected material

¹Rabin (1993, p. 1281), italics in the original.

payoff and a reciprocity component that is bounded above and below,² and which Rabin uses to model conditional altruism: A's reciprocity component is positive if A treats B kindly (unkindly) when she believes that B's intentions are kind (unkind). Further, A's reciprocity component collapses to zero if x_A^* is just equal to the equitable payoff.

Rabin (1993) resort to Psychological Game Theory -Geanakoplos et al. (1989)- to model the idea that beliefs about the other player's intentions affect utility, and proposes, in line with that theory, an equilibrium concept in which players' strategies are optimal given their beliefs which, moreover, turn out to be correct.³ Nevertheless, his solution concept is problematic in sequential games where non-optimizing behavior may be prescribed out of the equilibrium path. Dufwenberg and Kirchsteiger (2004) extend Rabin's approach to n-player extensive form games and provide a solution concept that follows the logic of subgame perfect equilibrium.

It follows from Rabin's definition of the equitable payoff, the reference point that players use to judge whether intentions are kind or not, that the whole set of allocations -including those outcomes that the other player does not intend to reach- might affect one's behavior, something that is generally compatible with experimental evidence. On the other hand, two assumptions of Rabin (1993) and Dufwenberg and Kirchsteiger (2004) are largely incompatible with experimental evidence. First, the equitable payoff is independent of the opponent's expected payoffs which implies, for instance, that if (2,1, 0) and (2, 2) are the only Pareto efficient material allocations and the second player's intentions are (2, 2) -i.e., he has unkind intentions towards the first player- then the first player might be willing to hurt the second one, if possible. A second shortcoming is that dummy players, which cannot have kind or unkind intentions, are never treated kindly (or unkindly). For instance, both models predict no giving in the dictator game.⁴

Falk and Fischbacher (2006) propose another extension of Rabin (1993) to extensive form games that avoids those two problems. In it, a player's utility at a terminal node z is the sum of the material payoff at z and the reciprocity utilities that she gets at all her decision nodes that precede z -players may weight such reciprocity utilities differently, thus introducing heterogeneity. The chooser's reciprocity utility at decision node n depends on her beliefs at n about the opponents' intentions. As in Rabin (1993), kind intentions

²Hence, the bigger the material payoffs, the less the players' behavior reflects their concern for fairness.

³Note well that, sensibly, beliefs are fixed exogenously and are not an object of choice.

⁴Nevertheless, appendix A of Rabin (1993) extends the main model to avoid this problem.

trigger *ceteris paribus* reward whereas unkind intentions trigger punishment. Contrary to Rabin (1993), however, intentions are kind (unkind) at n when the other player gets a lower (higher) expected material payoff than oneself's.

Another key distinction from Rabin's model is that the *intensity* with which agent A's wishes to punish or reward B depends on the whole set of outcomes that A believes that B believes to be reachable. Roughly, A's disposition to reward B lessens if A believes that, although B has kind intentions, B could have given more to A at *any* other available alternative. Conversely, A's disposition to punish B lessens if A believes that, although B has unkind intentions, B could *not* have given more to A at *any* other alternative without putting himself in a disadvantageous position -an "unreasonable sacrifice" by player B.

Models using Psychological Game Theory are based on the interesting idea that social emotions depend on the motives we attribute to others. However, they share the drawback of being rather complex. Levine (1998) improves tractability by assuming that people are concerned about the opponent's type and not about his intentions. A typical agent A's type is completely specified by a number $a_A \in (-1, 1)$, which signals whether someone is benevolent ($a_A > 0$) or malevolent ($a_B < 0$). Given this, and using the linear framework of equation (1), $\alpha_{AB}(z)$ depends positively on a_A and a_B . For example, even if player A is benevolent, she may become spiteful ($\alpha_{AB} < 0$) towards a sufficiently malevolent player B. Since the type of each player is private information, there is a possibility for signalling, that is, players' actions may reveal how benevolent (or malevolent) they are, and their opponents care about this. In that way, non-chosen moves may be as important as the moves one actually chooses, something that, as we have already remarked, is sensible and consistent with experimental evidence. One drawback of this model is that it renders a multiplicity of equilibrium *strategy profiles* in most games.

The appendix of Charness and Rabin (2002) offers another model of reciprocity. They introduce a demerit profile $d = (d_1, \dots, d_n)$, where $d_j \in [0, 1]$ for all j , and nonnegative parameters λ, δ, b, k, f where $\lambda \in [0, 1]$ and $\delta \in (0, 1)$. Player i 's utility function is

$$\begin{aligned}
 U_i = & (1 - \lambda) \cdot x_i + \lambda[\delta \cdot \min(x_i, \min_{k \neq i}\{x_k + bd_k\}) \\
 & + (1 - \delta)(x_i + \sum_{k \neq i} \max\{1 - kd_k, 0\} \cdot x_k) - f \sum_{k \neq i} d_k \cdot x_k].
 \end{aligned}$$

The key aspect of these preferences is that the greater is d_k for $k \neq i$, the less weight

player i places on player k 's material payoff. In fact, if f and d_k are sufficiently large then player i wishes to hurt player k .

In order to model reciprocity, Charness and Rabin endogenize each demerit d_j to make it dependent on player j 's strategy. Roughly speaking, they define $g_i(s_i, s_{-i}, d)$ as a correspondence selecting those values of $\lambda \in [0, 1]$ such that s_i is a best response to s_{-i} given demerits d . Each $g_i(s_i, s_{-i}, d)$ is then compared with an exogenous 'selflessness standard' λ^* -to be interpreted as the weight that a decent person *ought to* put on social welfare. The intuition is that if $\max\{g | g \in g_i(s_i, s_{-i}, d)\} < \lambda^*$ then other players resent player i 's choice. Given all this, strategy profile s is a 'reciprocal-fairness equilibrium' (RFE) if there exists a profile d and a correspondence $g_i(s, d)$ for all i such that, for all i , s_i is a best response to s_{-i} given d , and $d_i = \max[\lambda^* - g_i, 0]$ -i.e., the demerit profile must be consistent with the profile of strategies.

The model of Charness and Rabin (2002) presents several drawbacks like its complexity, the existence of many free parameters, the lack of heterogeneity in players' utility functions, and the fact that it is unclear how to compute utilities if there are multiple equilibrium demerit profiles. Charness and Rabin (2002, 851) do not see their model as "[...] being primarily useful in its current form for calibrating experimental data, but rather as providing progress in conceptualizing what we observe in experiments." In this respect, and because several of their intuitions are somehow present in our model, one may see it as a tractable continuation of their research.

All above mentioned utility models are *non-consequentialistic* or non-separable (Camerer, 2003) because a player's utility at terminal node z does *not* only depend on the distribution of material payoffs at z . Other models are consequentialistic or separable. Fehr and Schmidt (1999) and Bolton and Ockenfels (2000), for instance, model inequity aversion.

Two key hypothesis characterize Fehr and Schmidt (1999) -we use equation (1) to describe them. First, $\alpha_{AB}(z)$ is a *positive* parameter if A gets a *larger material* payoff than B , that is, if $x_A(z) > x_B(z)$. Second, $\alpha_{AB}(z)$ is a *negative* parameter if A gets a *smaller material* payoff than B -in other words, agents are envious. In addition, players are heterogeneous regarding the inequity aversion parameters, and for any individual, the envy parameter is larger than the parameter measuring advantageous inequity aversion.

Analogously, Bolton and Ockenfels (2000) posit two basic assumptions -we use again equation (1). Roughly speaking, $\alpha_{AB}(z)$ is positive (negative) if A 's material payoff is

above (below) the *average material* payoff at z . Note well that these two conditions hold independently of how big B 's material payoff is. Bolton and Ockenfels (2000) also assume that individuals are heterogeneous.

To finish, Andreoni and Miller (2002) and Charness and Rabin (2002) also offer consequentialistic models to explain evidence coming from some experiments. Charness and Rabin (2002), for example, hypothesize that $\alpha_{AB}(z)$ is positive and, moreover, $\alpha_{AB}(z) > \alpha_{AC}(z)$ for any other player C if B happens to be the worst off agent. That is, players are altruists with Rawlsian maximin concerns.⁵

3 The Model

Material Games and Norms. Consider any extensive form game of perfect recall. Let $N = \{1, \dots, n\}$ denote the set of players, and $u(z) = \{u_1(z), \dots, u_n(z)\}$ the vector of players' payoffs at terminal node z . Players are rational, that is, each one seeks to maximize her own payoff given her beliefs about other players' strategy.

In addition, let $x(z) = \{x_1(z), \dots, x_n(z)\}$ denote the vector of *material* payoffs at z . That is, $x_i(z)$ represents the cardinal utility that player i gets from consumption, money, and effort exerted along the history of z . Nevertheless, in lab games -our main concern here- it seems reasonable to simplify and assume that subjects' material welfare just coincides with earned money. Throughout the paper, hence, the terms 'monetary payment' and 'material payoff' are synonyms.

Material payoffs and payoffs are not the same thing -i.e., generally $x_i(z) \neq u_i(z)$ for any player i and node z . However, the researcher may initially have information only about the *material game*, that is, the researcher may know $x_i(z)$ for any i and z but not $u_i(z)$. We propose in what follows a theory on how to derive any $u_i(z)$ from the data contained in the material game.

Definition 1 *A norm Ψ is a nonempty correspondence $\Psi : h \rightarrow A(h)$ that applies on any information set h of any material game. Action $a \in A(h)$ is said to be consistent with norm Ψ if $a \in \Psi(h)$. Otherwise, a is a deviation from Ψ .*

⁵When mentioning Charness and Rabin (2002) in what follows, and unless otherwise noted, we refer to their model of quasi-maximin preferences and not to the previously described reciprocity model.

One may interpret a norm as a prescription indicating how one *ought* to behave at any conceivable situation at which one may be called to move. To put it like that, a norm orders the available actions at any information set: Some are commendable and others are not. We provide below an example of a specific norm: The Efficiency and Equity norm, or E-norm.

Preferences. To simplify matters, assume that the E-norm is the only norm in the society and that there exist two types of agents: Selfish and principled. Selfish people ignore the E-norm and just care for their material payoff. Therefore, the utility of any such player at node z is given by

$$u_i(z) = x_i(z).$$

On the contrary, principled people have internalized the E-norm and suffer a cost when violating it, to be interpreted as a painful emotion. Furthermore, the intensity of the emotion depends *inversely* on the number of transgressors. Thus, a principled deviator feels happier if every player deviates than if she is the only deviator. One can interpret these assumptions as modelling the effects of shame on preferences. In effect, in López-Pérez (2005) we provide psychological evidence and argue that a deviation from an internalized norm triggers shame and that shame intensity is strongly correlated with inferiority feelings -e.g., on how one's actions compare with others'.

To formalize this, let $R(z)$ designate the set of players that respected the norm in the history of z . Namely, $R(z)$ includes all players who made choices consistent with the norm *or* no choice at all in the history of z . Further, let $r(z)$ denote the cardinality of set $R(z)$. Given all this, a principled player's utility function takes the following form:

$$u_i(z) = \begin{cases} x_i(z) & \text{if } i \in R(z). \\ x_i(z) - \gamma \cdot r(z) & \text{if } i \notin R(z); (\gamma > 0). \end{cases}$$

Parameter γ measures how intensely principled types have internalized the norm. The larger it is, the more pain a principled deviator feels *ceteris paribus*. Importantly, γ is independent of the particular deviation oneself made in the past. Although this is indeed an extreme simplification, we show throughout the paper that it is enough to replicate many *qualitative* experimental results.

The E-norm. Let h denote an information set, $A(h)$ denote its corresponding set

of available actions, t_0 denote an *initial* decision node, that is, any node immediately succeeding Nature's moves -i.e., random shocks- and $X(t_0)$ denote the set of all $x(z)$ that succeed decision node t_0 .

Definition 2 Allocation $x = \{x_1, \dots, x_n\} \in X(t_0)$ is an (ε, δ) -**fairmax distribution** of a material game if it maximizes function

$$F_{\varepsilon\delta} = \varepsilon \cdot \sum_{i \in N} x_i - \delta (\max_{i \in N} x_i - \min_{i \in N} x_i), \quad (2)$$

over $X(t_0)$ for at least one node t_0 . A path connecting node t_0 and one of its (ε, δ) -fairmax distributions is an (ε, δ) -**fairmax path** of the material game.

Assuming $\varepsilon, \delta > 0$, function $F_{\varepsilon\delta}$ depends *positively* on the *social efficiency* of x - measured as the sum of monetary payoffs- and *negatively* on the degree of *inequality* embodied in x . In what follows, and given two real numbers a and b , F_{ab} designates function $F_{\varepsilon\delta}$ when $\varepsilon = a$, and $\delta = b$. Unless otherwise noted, we normalize the efficiency parameter ε to one and keep δ strictly positive but smaller than one. Assumption $1 > \delta$ indicates that social efficiency is relatively more important than equality (we argue later why this assumption *seems* to be reasonable). To simplify the exposition, we refer in what follows to a $(1, \delta)$ -fairmax distribution and a $(1, \delta)$ -fairmax path as a 'fairmax distribution' and a 'fairmax path', respectively.

To apply the E-norm to any material game, start by finding all its fairmax paths.⁶ Once this task has been completed, the E-norm selects actions as follows:

- (i) If information set h has *at least one* node on a fairmax path, the E-norm selects all actions of $A(h)$ that belong to a fairmax path.
- (ii) Otherwise, the E-norm selects the whole set $A(h)$.

⁶Infinite material games may have no fairmax distribution. Suppose, for example, that t_0 is such that $X(t_0)$ consists of all vectors (x_1, x_2) such that $x_1 + x_2 = 1$, *except* $x = (1/2, 1/2)$. It is trivial then that no distribution maximizes function $F_{1\delta}$ over $X(t)$ when $\delta \neq 0$. All the material games we consider in the applications have at least one fairmax distribution. For completeness, however, one may assume that the E-norm allows any move at any h of a material game with no fairmax distribution. For an alternative, consult López-Pérez (2005).

It is worthy to mention the ideas that are buried in this norm. First, it is a norm of distributive justice or fairness that sees fairness as positively depending on social efficiency and equality. Second, this norm commends any player to play fairly -i.e., to follow a fairmax path- *if others played fairly as well*. Of course, a player may be uncertain about previous play in some information sets of some games. To put it like that, the norm commends in such a case to put one's faith on any previous mover, and play *as if* one believed that every previous mover played fairly before. Finally, if the mover at h knows that *at least one* deviation has taken place then any action becomes commendable. This latter feature is indeed extreme but it is enough to get our results and simplifies the analysis. The model is flexible enough, however, to introduce other, more realistic norms. We discuss this point in a bit more of detail in the conclusion -see also López-Pérez (2005) for descriptions of alternative, more sophisticated norms.

Players' Information. Equilibrium Concept. We will often assume that each player's type is private knowledge. That is, prior to the start of any game Nature draws players independently from a population with a binomial distribution over the set of types. Let μ denote the *objective* probability of being a principled agent. Following standard usage, we assume throughout the main text that μ is common knowledge -i.e., priors are common. In the appendix, however, we explain one possible way to relax this assumption in order to introduce some more heterogeneity, together with some applications of this idea.

Unless otherwise noted, we will use Perfect Bayesian Equilibrium (PBE) as a solution concept. A PBE consists of a probability assessment (beliefs) over the nodes of each player's information sets and a strategy profile. Assessments reflect what the player moving at the corresponding information set believes has happened before reaching it. They must be, to the extent possible, consistent with Bayesian updating on the hypothesis that the equilibrium strategies have been used to date. In addition, any player's strategy in a PBE must be sequentially rational. That is, everybody must choose optimally at any of her information sets given her beliefs at that set and the fact that future play will be governed by the equilibrium strategies.

To finish, it is important to clarify that in our model assessments do not play any practical role at information sets *out of a fairmax path*. That is, once a "bad" action has taken place and that becomes common knowledge, beliefs about the type of the opponent

are unimportant to explain behavior. Because of this, we will not mention such beliefs when describing a PBE. This simplifies considerably the analysis.

4 Explaining Experimental Evidence (I)

In this section we use the model to explain experimental evidence coming from a number of games. In addition, we will provide some tentative answers to three questions: (1) How well does the E-norm approximate the actual moral standards that some subjects take to the lab? (2) What are the factors that explain norm compliance in one-shot games if deviations cannot be punished? and (3) Do people treat equally well passive players and active *and compliant* ones?

4.1 On Individual Decision Problems: Efficiency and Equality Matter

Why have we assumed that the E-norm is the only norm in the society? One reason is indeed parsimony: We could assume that principled types are heterogeneous and care about different norms, but this would complicate the model. A second reason is that the available empirical evidence *seems* to fit rather well with the idea that social efficiency and equality are important ingredients of subjects' views on distributive justice, and that social efficiency is relatively more important than equality.

To justify this last statement, it is convenient to consider the simplest possible scenario: An individual decision problem with externalities. In general terms, our model predicts that selfish agents choose the allocation that maximizes their own monetary earnings whereas principled agents choose the same allocation that selfish ones if their guilt parameter γ is sufficiently low and the fairmax distribution *that gives them a highest monetary payoff* otherwise.

An example will help to understand these predictions. In the so-called Dictator game, one subject (the 'dictator') is endowed with a sum of money M and must decide how much of that money to transfer to another subject (the 'dummy'). Obviously, selfish dictators give nothing. What about principled ones? As the unique fairmax distribution is equal sharing, a principled dictator will get a utility payoff of $\frac{M}{2}$ if she gives half of the money to the other subject, and a utility payoff of $M - x - \gamma$ if she transfers $x \neq \frac{M}{2}$ (to understand

why she suffers a psychological cost γ , note that the dummy belongs to set $R(z)$ for any z because she makes *no* choice in the game). This latter expression takes a maximum value of $M - \gamma$ when $x = 0$, and hence it follows that a principled dictator respects the E-norm if the psychological cost γ is larger than $\frac{M}{2}$, gives nothing if γ is smaller than $\frac{M}{2}$, and is indifferent between both options if γ equals $\frac{M}{2}$. To sum up, principled dictators follow their principles if that is not too costly.⁷

Experimental results on the Dictator game -see Camerer (2003), pp.57-58, for an extensive survey- are somewhat sensitive to the degree of anonymity enjoyed by subjects when choosing, and the wording of instructions. Nonetheless, one may reasonably contend that (i) the average offer is around $0.25M$, (ii) an average of 35-40% of the participants give nothing, and (iii) there are virtually no offers above 50% of the stake. Result (iii) is clearly replicated by our model, whereas results (i) and (ii) are consistent if we assume that μ and γ take appropriate values.⁸

Our predictions depend heavily on the values of the efficiency parameter ε and the inequality parameter δ of function (2). To illustrate this, assume for a moment $\varepsilon = 1$ and $\delta = 0$. Since any monetary allocation in the dictator game is $(1, 0)$ -fairmax, the model would then forecast that *any* type of player gives zero money. Dictator game results, therefore, reject a model based on the idea that *all* principled players believe that distributive justice *exclusively* depends on social efficiency. If, on the contrary, one assumes $\varepsilon = 0$ and $\delta = 1$ (that is, equality is the only ingredient of justice), equal sharing is the only $(0, 1)$ -fairmax distribution and predictions coincide with those when $\varepsilon = 1$ and $0 < \delta < 1$.

Therefore, the standard dictator game does not discriminate between a formulation based on function F_{01} and the one we use throughout the paper, based on $F_{1\delta}$ for $0 < \delta < 1$. In contrast, the ingenious design of Andreoni and Miller (2002) allows for that. In their dictator game experiments, transfers of money were multiplied by a factor that differed from session to session and was common knowledge. In one session, for instance, the

⁷Hence, marginal changes in parameter γ may produce radical switches in principled agents' behavior. This feature disappears if we assume that the intensity of the internal punishment conveniently depends on the particular deviation a principled agent does. We have investigated this issue in López-Pérez (2005).

⁸However, our model fails to provide an accurate picture of the actual distribution of offers, which are usually scattered along the interval $[0, M/2]$ and not concentrated on the extremes. See the previous footnote to this respect.

factor was equal to 3 so that a transfer of x units of the dictator’s initial endowment translated in final earnings of $3x$ for the dummy. In this session and also when the factor was equal to 2, a significant number of dictators made transfers such that they ended up with less money than the receiver. This is consistent with our specification based on $F_{1\delta}$ for $0 < \delta < 1$ but not with one based on F_{01} . Thus, agents seem to be concerned with *both* social efficiency and equality, assigning a larger weight to the first variable.

More experimental data supports this conjecture. In Study 2, Decision 1 of Charness and Grosskopf (2001), subjects had to choose between (self, other) allocations of pesetas (625, 625) and (600, 1200).⁹ Trivially, selfish agents choose allocation (625, 625) whereas principled ones choose the efficient allocation (600, 1200) *if* γ and δ are high and small enough, respectively, and the egalitarian one otherwise. Charness and Grosskopf (2001) report that only 33.3% of the subjects (N=108) chose the egalitarian allocation.

Additionally, in their Study 2, Decision 3, the *same* subjects received 600 pesetas and had to choose any payoff for another participant between 300 and 1200 pesetas. 74.1% of the subjects chose 1200 pesetas -i.e., the only fairmax distribution. Only 10.2% of them chose opponent’s earnings equal to 600 pesetas -i.e., the egalitarian distribution.

To sum up: *If one assumes that principled types exclusively care about one norm*, a model based on our E-norm is more coherent with the data than one based on a norm of social efficiency (F_{10}), a norm of equality (F_{01}), or on a norm of social efficiency and equality in which equality receives a larger weight than efficiency (that is, $F_{1\delta}$ for $1 < \delta$). It is important to stress, however, that there exist other norms that we have not considered here and that *might* generate better predictions. Suppose, for instance, that we had defined a fairmax distribution as an allocation maximizing function

$$Q(x) = \sum_{i \in N} x_i + \tau \min_{i \in N} \{x_i\}, \quad (3)$$

where $0 < \tau$. Closest to Charness and Rabin’s theoretical approach, this function combines social efficiency and maximin. Would such a model explain better the data? Note first that in two-player games this formulation predicts similar results to one based on our E-norm. Therefore, we need multiple-player games to discriminate between both

⁹At the exchange rate of the moment, one US dollar was around 150 pesetas. Each participant took decisions in three different problems but was paid for only one of those problems, chosen at random at the end of the session.

formulations. For instance, consider four players A, B, C, D and assume that A is the only active player, that she is principled, and that she must choose between (A, B, C, D) dollar allocations (100, 131, 130, 49) and (100, 200, 50, 50). Both allocations are equally efficient, but the first one is more egalitarian (at least according to function $F_{1\delta}$), while the second one maximizes the payment of the worst-off player. As a result, A would choose the former allocation if she had internalized our E-norm, and the latter one if she cared about a norm based on the above mentioned function $Q(x)$ (note that she gets 100 dollars whatever her choice so that her value of γ is immaterial here). Unfortunately, and as far as we are concerned, there is no empirical data on dictator games like this one that would allow us to discriminate between these two approaches.¹⁰ Hence, this is still an open question that requires more attention from experimental economists.

We finish by comparing our predictions with those from other models. For instance, Fehr and Schmidt (1999), Bolton and Ockenfels (2000) and Falk and Fischbacher (2006) assume no efficiency concerns and a more or less complex form of inequity aversion. It should be hence clear that they are inconsistent with the evidence cited previously. The same occurs with Levine (1998), at least if we take the distribution of types that Levine posits. Rabin (1993) and Dufwenberg and Kirchsteiger (2004) are also inconsistent with the above mentioned data because they predict that no agent sacrifices her own material payoff to reward a dummy player -this is true at least for the simplest version of Rabin (1993). On the contrary, Charness and Rabin (2002) is largely consistent.

Charness and Rabin (2002) also report abundant experimental evidence that contradicts a utility model based *exclusively* on (linear) inequity aversion and material interest. In game Berk23 of Charness and Rabin (2002), for instance, subjects choose between (self, other) allocations of US dollars (2, 8) and (0, 0). Our model predicts that *all* agents choose the first allocation and this is exactly the actual result. In contrast, Fehr and Schmidt (1999), and Bolton and Ockenfels (2000) predict that a significant proportion of subjects choose (0, 0). Note further that Levine (1998), which assumes that some players are benevolent while others are malevolent, shares this prediction.

As another illustration, participants in game Barc2 of Charness and Rabin (2002) chose between (self, other) allocations of pesetas (400, 400) and (375, 750). Principled

¹⁰For instance, Engelman and Strobel (2004) study dictator games with multiple dummies, but choices are not efficiency-preserving in any of their games.

players choose the second allocation if γ is high enough whereas selfish agents choose the first allocation. On the opposite, inequity aversion models predict that *all* agents choose the first allocation. It turned out that 50% of the participants chose the first allocation. In the same line, Charness and Rabin (2002) show that 69% of the participants choose (Self, Other) allocations of dollars of (4, 7.5,) over (4, 4).

4.2 Norm compliance without punishment threats

People respect norms in nonrepeated interactions even if transgressions cannot be punished and compliance is contrary to material interest. This section explores what factors affect norm compliance in such settings. Intuitively, a principled player will obey the E-norm -and, by extension, any other norm- if two conditions hold.

First, she must have internalized the E-norm with enough intensity. In effect, principled agents suffer a psychological cost if they deviate from the E-norm. Nevertheless, if the expected material benefit of deviating is sufficiently high to overcome the expected pang - which depends, among other factors, on parameter γ - she may succumb to the temptation and deviate. Consequently norm compliance requires sufficiently *strong convictions*.

Second, she must believe that sufficiently many other players will comply as well. This follows from the fact that the psychological cost of deviating depends directly on the number of norm followers, and means that norm compliance follows a *reciprocal* logic. Let us also remark that a player's expectations that the others will comply subtly depend on her expectations about the other players' *types*. Believing that player B is selfish suffices to infer that B will indeed deviate. On the contrary, believing that B is principled is *not* enough to sustain the belief that B will comply. We will come back to this point later.

To illustrate all those points, consider a two-player material game in which players 1 and 2 must choose simultaneously positive numbers q_1 and q_2 , respectively. As a result, player i gets a monetary reward $x_i = Kq_i - q_iq_j - q_i^2$, where K is a positive number and $i, j \in \{1, 2\}, i \neq j$. This is a Cournot duopoly game in which firms' marginal costs take a common, constant value c and demand is linear, that is, $p = M - Q$, where M is a constant ($M > c$), p denotes the price, and Q the sum of quantities produced by each firm. In this setting, $K = M - c$.

Standard optimization techniques show that the sum of monetary rewards is maximized when $q_1 + q_2 = \frac{K}{2}$. Moreover, each player gets the same material payoff if $q_1 = q_2$.

Therefore the unique fairmax distribution of this game is implemented when both players choose $q_F = \frac{K}{4}$, and that is what the E-norm commends.

When do players respect the E-norm? A first point to make in this respect is that no player will do that if she does not expect the opponent to comply as well. Basically, this occurs because producing q_F is never in the firms' material interest -a standard textbook result shows that *if the firm is selfish* then q_F is a strictly dominated strategy. Consequently, selfish firms will never produce q_F , and principled ones will choose q_F only if they expect the opponent to produce q_F as well.

Proposition 1 *A strategy profile in which both selfish and principled players choose $q_{NC} = \frac{K}{3}$ is a PBE for any prior μ . This is the only PBE strategy profile in which both types deviate from $q_F = \frac{K}{4}$.*

Proof. Selfish agents always seek to maximize material reward $x_i = Kq_i - q_iq_j - q_i^2$. The same is true for principled agents if the opponent deviates from the norm ($q_j \neq \frac{K}{4}$). Fixing q_j , one may show by standard optimization techniques that maximization of x_i requires $q_i = \frac{K-q_j}{2}$. Now, given the symmetry of the problem, both players should make the same choice at equilibrium. Hence, we have $q_1 = q_2 = \frac{K}{3}$. ■

Production level q_{NC} corresponds to the textbook *Nash-Cournot* prediction when both firms are self-interested. Further, note that this equilibrium exists for any μ and, in particular, for $\mu = 1$, that is, in a *complete* information game played by two principled agents. Intuitively, and since principled types follow norms reciprocally, mutual *distrust* -i.e., the mutual expectation that the opponent will not comply- destroys any respect for the norm. Assuming then that principled types trust each other, does an equilibrium exist?

Proposition 2 *A strategy profile in which principled players choose $q_F = \frac{K}{4}$ and selfish ones $q_C = \theta q_F$, where $\theta = \frac{4-\mu}{3-\mu}$, is a PBE strategy profile if $4(3-\mu)\sqrt{\mu\gamma} \geq K$.*

Proof. We show first that selfish agents play a best-response. Expected utility of playing q_i is given by

$$\mu(Kq_i - q_iq_F - q_i^2) + (1-\mu)(Kq_i - q_iq_C - q_i^2) \quad (4)$$

or alternatively by $q_i\pi - q_i^2$, where $\pi = K - \mu q_F - (1-\mu)q_C$. Differentiating $q_i\pi - q_i^2$ with respect to q_i and equating that to zero, we get as a necessary (and sufficient) condition for maximum that $q_i = \frac{\pi}{2}$, and some algebra shows that $\frac{\pi}{2} = q_C$.

To prove that principled agents play a best response as well, we first compute their expected utility of playing q_F . For that, and since they follow the norm and feel no remorse, it suffices to substitute q_F for q_i at expression (4) so that one gets that expected payoff equals $\pi q_F - q_F^2 = q_F^2(2\theta - 1)$.

Suppose now that a principled agent deviates from q_F . Her expected utility is then $q_i\pi - q_i^2 - \mu\gamma$ and the best she can do is producing $q_i = \frac{\pi}{2}$, hence getting an expected payoff of $\frac{\pi^2}{4} - \mu\gamma = \theta^2 q_F^2 - \mu\gamma$. It follows that playing q_F is optimal if

$$q_F^2(2\theta - 1) \geq \theta^2 q_F^2 - \mu\gamma. \quad (5)$$

And some algebra proves inequality (5) to be equivalent to $4(3 - \mu)\sqrt{\mu\gamma} \geq K$. ■

Note that this equilibrium exists only if parameters γ and μ are large enough. Experimental evidence on the Cournot game is summarized in Holt (1995). Although results are far from conclusive, they show that a significant number of participants in one-shot games attempt tacitly to collude, choosing output levels close to the joint-income maximizing level q_F , whereas remaining subjects make quantity choices around the Nash-Cournot equilibrium. Interestingly, if repetition (with rematching) is allowed, cooperation tend to vanish with time and most output decisions shift back to the Cournot level. However, the rate of convergence is not equal for all pairs of subjects: Some of them converge very fast while others mutually cooperate for a significant number of rounds. We briefly argue in the appendix that this might be explained by belief heterogeneity and learning about the opponent's type.

4.3 Positive reciprocity: Active and Passive Players

Most modern models of reciprocity -Rabin (1993), Levine (1998), Dufwenberg and Kirschsteiger (2004), and Falk and Fischbacher (2006)- predict a phenomenon called *positive reciprocity*. That is, people are in average more generous and kind towards those who exhibited kind behavior in the past than towards passive players that did not perform any action -the reciprocity model of Charness and Rabin (2002) is an exception because if a player does not misbehave then her demerit is zero, as a dummy player's. The E-norm, on the contrary, only allows "unkind behavior" if it is certain that the opponent deviated before. This implies that dummy players are as equally legitimated to receive a kind treatment as active norm compliers -i.e., our model predicts *no* positive reciprocity.

To illustrate the differences between our model and other reciprocity models, consider the mini-trust *material game* represented at Figure 1. The first mover (the ‘investor’) chooses either not to trust (move D) or to trust (move T) the second player (the ‘trustee’). In the first case, the investor gets x monetary units and the trustee gets 0 units. Alternatively, the investor may trust and give the trustee the chance to repay trust (move R) or not (move A). If trust is repaid, both earn r ($> x$) monetary units. If trust is not repaid, the investor gets the ‘sucker’ payoff s ($< x$) and the trustee earns the highest payment t . To sum up, we have $s < x < r < t$. In most experiments, values are chosen so that (r, r) is the unique fairmax distribution for any $\delta < 1$. We assume that in what follows.

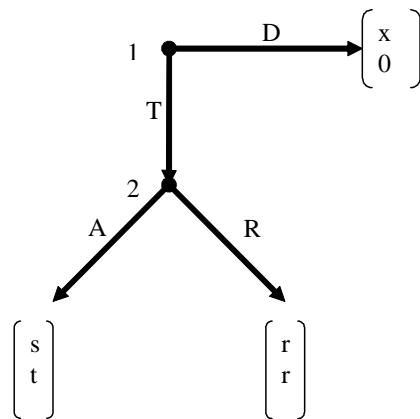


Figure 1. Mini-Trust Material Game.

Consider now two variations of this mini-trust game. In the *intentions* treatment player 1 is active and she effectively chooses her move whereas in the *random* treatment player 1 is passive and her move is decided by Nature -e.g., with the flip of a coin. Note that the unique fairmax path of the random treatment simply consists of action R -recall that a fairmax path always starts after all random moves have been made- whereas the only fairmax path of the intentions treatment is formed by moves T and R. Suppose then that player 2 is asked to move, will he behave differently in *each treatment*? The answer is negative.

Proposition 3 *In equilibrium and independently of the treatment, selfish trustees choose A whereas principled ones choose R if γ is high enough and A otherwise.*

Proof. Selfish movers go for the highest material payoff so that they play A. Since the E-norm commends trustees to move R in both treatments, principled trustees comply if the utility of playing R is larger than that of playing A, that is, $r > t - \gamma$. ■

Let us compare with other models. Consequentialistic models as Fehr and Schmidt (1999), Bolton and Ockenfels (2000) and Charness and Rabin (2002) predict no behavioral difference between both treatments. On the contrary and as we stated above, most models of reciprocity predict a significant decay in repay in the random treatment. Levine (1998), for instance, predicts some decay because trusting in the intentions treatment signals benevolence, which is rewarded, whereas the random treatment does not allow this kind of type-selection.

Most of the evidence in this regard is consistent with our model. Dufwenberg and Gneezy (2000) report data from an experimental Lost Wallet game which is very similar to our mini-trust game -the main differences are that player 2 faces a continuum of choices (more precisely, he plays a dictator game with stake size $2r$) if player 1 trusts him, and that x is larger than r in some treatments (however, $x < 2r$), but these differences are inconsequential for our model. They compare second movers' choices with data from a pure dictator game with stake size $2r$ and do *not* reject the hypothesis that both sets of data come from the same distribution. Since a pure dictator game can be seen as a particular case of our random treatment, this is plainly consistent with proposition 3. Dufwenberg and Gneezy (2000) also show that second movers' payback is uncorrelated with player 1's outside option x , and this is again consistent with our model but not with other reciprocity models -except that of Charness and Rabin (2002). Charness and Rabin (2002), Offerman (2002), and Cox and Deck (2005) report similar results. In contrast, a number of experimental papers have reported opposite results -consult Camerer (2003, pp. 89-90) for a useful discussion and references.¹¹ More experimental research is due in this regard.

4.4 Games with multiple moves: The Centipede Material Game

This is a two-person material game that resembles the homonym game introduced by Rosenthal (1982). For expositional purposes, assume that player 1 (2) is female (male). Each player alternately gets a turn to either terminate the game or pass the turn to the opponent. In the original version, this process may last a maximum of one hundred moves (hence the name "centipede"). In the last node, player 2 chooses between a socially

¹¹This evidence suggest that people care *sometimes* about others' intentions, but also that this concern seems to be rather fragile.

efficient monetary vector (x_1, x_2) and an inefficient one (x_1^*, x_2^*) in which, however, he gets a larger monetary payment, that is, $x_1 + x_2 > x_1^* + x_2^*$ but $x_2^* > x_2$. In turn, if player 1 terminates the game at the penultimate node, distribution (x'_1, x'_2) is reached such that $x_1^* + x_2^* > x'_1 + x'_2$ but $x'_1 > x_1^*$. An analogous pattern repeats in all the previous nodes.

If it is common knowledge that both players are selfish and rational, there exists a unique subgame perfect equilibrium. In it, any player terminates at any node and, consequently, the least efficient outcome is reached. One can show this by backwards induction. In effect, a selfish second player would choose the inefficient but own-payoff maximizing allocation at the last node. In turn, the first player would terminate the game in the previous node in order to get a higher payoff. Analogously, any player would choose, if given the choice, to terminate the game at any node.

Nonetheless, experiments with a simpler version of the centipede game show that this gloomy prediction is far from correct. McKelvey and Palfrey (1992) report experimental results from four move and six move centipede *material* games -Figure 2 reproduces the decision tree and *monetary* payments of the four move game. Most subjects choose to pass the turn in the initial nodes. Furthermore, a nonnegligible proportion of the participants pass at *every* decision node -including the last one- if given the opportunity, thus providing clear evidence that some subjects are not purely selfish. Finally, the proportion of movers who terminate the game at each node increases as the last node gets closer. Note that these results are also inconsistent with models of inequity averse preferences like Fehr and Schmidt (1999), and Bolton and Ockenfels (2000) which predict, at this specific centipede game, that all players terminate at the first node.

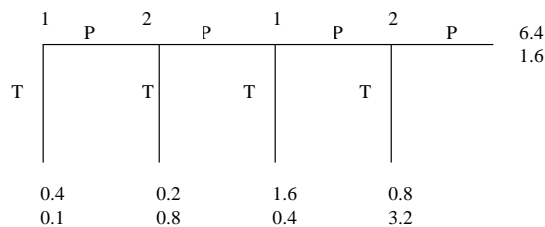


Figure 2. Four move Centipede Material Game.

We prove in what follows that if some subjects have internalized the E-norm, then the model can replicate the previous experimental facts and, moreover, we highlight that

players' degree of internalization of the norm -measured by parameter γ - and the priors μ are the key variables that explain individual behavior. To start, notice that there is only one strategy profile which is consistent with the E-norm: All players should pass the turn at every node.

Some required notation follows. We denote the initial node as node 1, the next node as node 2 and so on. The action of passing (terminating) at node k is denoted as P_k (T_k). Finally, let $\mu(k)$ denote the beliefs that the player moving at node k has about the opponent's type -obviously, $\mu(1) = \mu$. I study two cases: $0.8 < \gamma < 1.6$, and $1.6 < \gamma$ in order to show how differing levels of internalization affect behavior. The results can be extended to other values of γ .

Proposition 4 *Assume $0.8 < \gamma < 1.6$. If $\mu > 1/7$ there is only one PBE strategy profile. In it, player 2 plays (P_2, T_4) independently of his type, and player 1 plays (P_1, P_3) if she is principled and (P_1, T_3) if she is selfish. If $\mu < 1/7$ there is also a unique PBE strategy profile. Principled players play the same strategy as before. A selfish second mover plays T_4 and randomizes between T_2 and P_2 , assigning probability $\frac{1-7\mu}{7[1-\mu]}$ to T_2 . In turn, a selfish first player chooses T_3 and randomizes between T_1 and P_1 , assigning probability $\frac{6\mu}{1-\mu}$ to P_1 . In the marginal case in which $\mu = 1/7$, there exist multiple PBE.*

Proof. Since $\gamma < 1.6$, principled players pass at every node except the last one independently of their beliefs. Trivially, a selfish second player also moves T_4 and, consequently, a selfish first mover chooses T_3 . All this implies that a selfish second player moves T_2 if

$$0.8 > 3.2\mu(2) + [1 - \mu(2)](0.4) \Leftrightarrow \mu(2) < \frac{1}{7},$$

he moves P_2 if $\mu(2) > \frac{1}{7}$, and he is indifferent between T_2 and P_2 if $\mu(2) = \frac{1}{7}$. Moving backwards to the first node, we must consider three cases depending on the value of $\mu(2)$.

First, if $\mu(2) > \frac{1}{7}$ then a selfish first mover moves P_1 because any type of second mover plays P_2 subsequently. Since principled types also play P_1 , Bayes' law implies $\mu(2) = \mu$ so that consistency of beliefs requires $\mu > \frac{1}{7}$.

Second, if $\mu(2) < \frac{1}{7}$ then a selfish second player moves T_2 . Consequently, a selfish first mover plays T_1 if $0.4 > 1.6\mu + [1 - \mu](0.2) \Leftrightarrow \mu < \frac{1}{7}$, she plays P_1 if $\mu > \frac{1}{7}$, and she is indifferent between both actions if $\mu = \frac{1}{7}$. None of them can happen in equilibrium if $\mu(2) < \frac{1}{7}$. In effect, suppose first $\mu < \frac{1}{7}$. Then a selfish first mover plays T_1 so that Bayes'

law entails $\mu(2) = 1$, hence contradicting our hypothesis that $\mu(2) < \frac{1}{7}$. An analogous line of reasoning applies to remaining cases.

Third, assume $\mu(2) = \frac{1}{7}$ so that a selfish second mover is indifferent between P_2 and T_2 , and let $\rho(2)$ denote the probability that he chooses P_2 . A selfish first player chooses P_1 if

$$0.4 < 1.6\mu + [1 - \mu](1.6\rho(2) + [1 - \rho(2)]0.2) \Leftrightarrow \rho(2) > \frac{1 - 7\mu}{7[1 - \mu]}.$$

In that case, Bayes' law requires $\mu(2) = \mu = \frac{1}{7}$ and any strictly positive value of $\rho(2)$ is then optimal. In turn, a selfish first mover would choose T_1 if $\rho(2) < \frac{1-7\mu}{7[1-\mu]}$, and Bayes' law would then imply $\mu(2) = 1$, contradicting our assumption. Finally, if $\rho(2) = \frac{1-7\mu}{7[1-\mu]}$ then a selfish first mover is indifferent between T_1 and P_1 -let $\rho(1)$ denote the probability that a selfish player 1 chooses P_1 . Bayes' law implies

$$\mu(2) = \frac{\mu}{\mu + (1 - \mu)\rho(1)},$$

which gives $\rho(1) = \frac{6\mu}{1-\mu}$ in equilibrium since $\mu(2) = \frac{1}{7}$ by assumption. Further, note that consistency of beliefs requires $\mu < \frac{1}{7}$, thus concluding the proof. ■

The intuition why *principled people* pass in all nodes node but the last one should be clear. Since the amount of money at play in the first three nodes is so small compared with parameter γ , they pass the turn -i.e., respect the norm- to prevent feeling badly for being the deviator. In the last node, on the contrary, the material temptation to deviate is large enough to compensate the psychological cost they suffer as a result.

All this has an interesting implication: If the game grows in length, it requires stronger internalization of the norm -i.e., larger γ - to make principled agents pass the turn at the final nodes. In other words, if the size of the 'cake' grows very much then principled people need to be 'very principled' not to deviate at the final nodes. In fact, experimental evidence seems to show that people are not so strongly 'principled'. McKelvey and Palfrey (1992) compare behavior in a four and six move centipede game -the latter one obtains from the former (figure 2) by adding two additional nodes so that T_5 , T_6 , and P_6 lead to pecuniary allocations (6.4, 1.6), (3.2, 12.8), and (25.6, 6.4), respectively. According to their results, "If we compare the four move games to the *last* four moves of the six move games, there is more taking in the six move games."¹²

¹²McKelvey and Palfrey (1992), p. 810. They propose an incomplete information model in which some players are altruists that pass at every node whereas the remaining agents are selfish ones. Their model

Our model predicts the same phenomenon if, instead of adding additional turns, one simply multiplies all monetary payments by a common factor larger than one. For example, take the game at Figure 2 and multiply all its payments by 4. Parameter γ should be then larger than 6.4 for a principled type to pass at the last node. In comparison, γ only needs to be larger than 1.6 for that to happen in the game at figure 2. Of course, something similar is true for every node, not only the last one. Therefore, the probability of terminating at any node is theoretically equal or larger for the transformed game than for the game at figure 2. The experimental results that McKelvey and Palfrey (1992) offer are on this line, although the difference is not completely significant.

The reason why selfish agents pass the turn in the initial nodes follows from a reputational argument *a la* Kreps et al. (1982). As they know that principled players are willing to pass, they mimic such behavior in order to induce in a potential selfish opponent the belief that they are also principled people. In such a way, they succeed in making a selfish opponent pass and, hence, earn more money. Consistent with experimental data, the probability that a selfish agent abandons her mimicking strategy and terminates increases as she approaches the last node.

Proposition 5 *Assume $1.6 < \gamma$. In any PBE, a principled player passes at every node. If $\mu > \frac{1}{7}$ then a selfish second mover plays (P_2, T_4) whereas a selfish 1 plays (P_1, P_3) . If $\frac{1}{49} < \mu < \frac{1}{7}$, then a selfish second player moves T_4 and randomizes between T_2 and P_2 , assigning probability $\frac{1-7\mu}{1-\mu}$ to T_2 . She plays P_1 and randomizes between T_3 and P_3 , assigning probability $\frac{6}{7(1-\mu)}$ to T_3 . If $\mu < \frac{1}{49}$, a selfish first player assigns some probability to T_1 . Apart of that, the PBE strategy profile is identical to that when $\frac{1}{49} < \mu < \frac{1}{7}$. In the marginal cases in which $\mu = \frac{1}{49}$, and $\mu = 1/7$, there exist multiple PBE.*

Proof. The proof of this proposition is very similar to that of Proposition 4. It is left to the reader. ■

The previous proposition shows how the intensity of internalization, measured by parameter γ , affects behavior of both principled and selfish players. Another factor that obviously runs the results is the particular norm we posited. If one assumed a pure norm of equality (F_{01}), for instance, predictions would then coincide with the standard

has similar predictions to ours, although it cannot explain why the probability of terminating at the last node is lower in the six move game than in the four move game.

equilibrium. To understand why, observe that such norm recommends any player to terminate the game at any node. As a result, material interest would lead to the standard prediction. The fact that subjects do not behave like that suggests that they understand justice or fairness as something more than pure egalitarianism.

5 Extending the Model: Anger and Punishment

We assume that only principled agents -i.e., those who have internalized the E-norm- display anger. This hypothesis is somewhat speculative, although there is some supporting evidence coming from Burnham (1999). In this experiment, subjects played a constrained ultimatum game where the only two offers were either \$5 or \$25 out of \$40. Moreover, subjects' testosterone levels were measured using saliva samples. Now, it is well known that high levels of testosterone are correlated with aggressive behavior. Therefore it is not surprising that subjects with high testosterone levels were relatively more likely to reject the \$5 offer. But Burnham (1999) shows something more: Subjects with high levels of testosterone were also relatively more likely to make an offer of \$25. This kind of correlations are consistent with our model.

Hence, one only needs to introduce some changes in principled agents' utility function, which is now given by

$$U_i(z) = \begin{cases} x_i(z) & \text{if } R(z) \equiv N \\ x_i(z) - \gamma \cdot r(z) & \text{if } i \notin R(z); (\gamma > 0) \\ x_i(z) - \alpha \max_{j \notin R(z)} x_j(z) & \text{if } R(z) \subset N, i \in R(z); (1 \geq \alpha > 0). \end{cases}$$

Since anger goes associated with a desire to punish the deviator, we model it as history-dependent spite. Clearly, parameter α measures aggressiveness. Further, and for simplicity, anger intensity does *not* depend on the particular deviation that the deviator made, and angry agents focus at the best off deviator. These assumptions are probably a bit unrealistic, but that does not prevent the model from explaining much qualitative evidence -we will briefly discuss later how to extend the model. We maintain the remaining assumptions of the model that were introduced in the previous chapter.¹³

¹³One may wonder whether previous results, obtained without the anger assumption, still hold. With a small caveat, the answer is positive for two reasons. First, if a deviation from the E-norm occurs in any of the games we studied then the action that maximizes the material payoff of the nondeviator also minimizes the deviator's material payoff. Hence, an angry player would make the same choices as a selfish

5.1 Determinants of Punishment: The Ultimatum Game

In this sequential material game player 1 (the ‘proposer’) is provisionally allocated $M > 0$ monetary units and has to propose how to divide that money between her and player 2 (the ‘responder’). Given any proposal of sharing $(x_1, M - x_1)$, the responder can either accept or reject. If he accepts, he gets $M - x_1$ and the proposer gets x_1 . If he rejects, both get nothing.

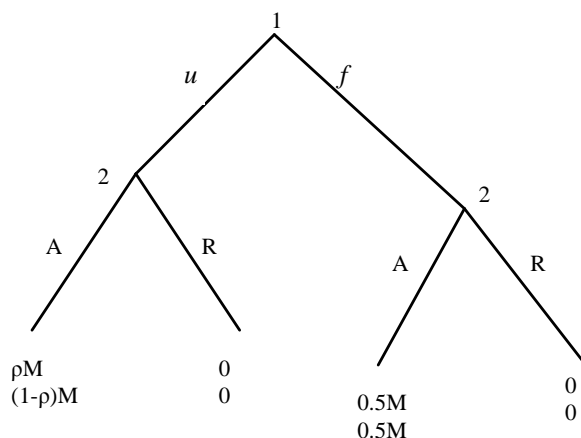


Figure 3. Mini-ultimatum material game.

Having a continuum of offers does not change our predictions. Thus, we have represented at Figure 3 a reduced version of the ultimatum material game in which player 1 has only two choices available: ‘Unfair’ (u) and ‘fair’ (f). The former choice consists of an offer of $(1 - \rho)M$ monetary units to player 2, where ρ is a number in the interval $[0, 1]$. In turn, choice f consists of a proposal of equal sharing and it is thus consistent with the E-norm. Player 2 can either accept (A) or reject (R) -note that if player 1 offers u then the E-norm allows player 2 to choose both A and R, whereas if player 1 chooses f the E-norm only allows acceptance.

Proposition 6 *For any priors μ and almost any ρ , the mini-ultimatum game has a unique PBE strategy profile. A principled responder always accepts the equal sharing*

one -note that this is no longer true for the games that we study in what follows. Second, if no deviation has taken place there is no place for anger, except as an expected emotion. Now, if someone expects the opponent to deviate then she will be less willing to respect the norm, because she expects to be angry, which is painful. To respect the norm, therefore, warm people require a larger parameter γ or a larger prior μ . Except for this quantitative differences, previous results still hold.

whereas he accepts an offer of $(1 - \rho)M$ if $\rho < \frac{1}{1+\alpha}$ and rejects if $\rho > \frac{1}{1+\alpha}$. A selfish responder accepts any offer if $\rho > 0$. A principled proposer's choice depends on the values of γ and ρ :

$\rho < \frac{1}{2}$ then she offers the equal sharing.

$\frac{1}{2} < \rho < \frac{1}{1+\alpha}$ then she offers u if $\gamma < \frac{M(2\rho-1)}{2}$ and f otherwise.

$\frac{1}{1+\alpha} < \rho$ then she offers u if $\gamma < \frac{M[(1-\mu)2\rho-1]}{2}$ and f otherwise.

Finally, a selfish proposer offers u if ρ and μ are large and small enough, respectively, and the equal sharing otherwise.

Proof. Accepting offer f is consistent with the E-norm and beneficial in monetary terms. Therefore, any type of responder accepts that offer. Offering u makes a principled responder angry so that he accepts u only if

$$(1 - \rho)M - \alpha\rho M > 0.$$

Trivially, a selfish responder accepts u for any $\rho > 0$.

Consider now a selfish proposer. She offers f if $\rho \leq \frac{1}{2}$ because f is always accepted and gives more money in that case. For analogous reasons, she offers u if $\frac{1}{1+\alpha} > \rho > \frac{1}{2}$. Further, offer u is not accepted by a principled responder if $\rho > \frac{1}{1+\alpha}$ so that a selfish proposer will make that offer only if $(1 - \mu)\rho M > 0.5M$, that is, if $\frac{2\rho-1}{2\rho} > \mu$.

To finish, consider a principled proposer. The 50-50 offer is clearly optimal if $\rho \leq \frac{1}{2}$. Finally, offering u gives $\rho M - \gamma$ units of utility if $\frac{1}{1+\alpha} > \rho > \frac{1}{2}$ and $(1 - \mu)\rho M - \gamma$ units of expected utility if $\rho > \frac{1}{1+\alpha}$. Simple algebraic manipulations prove that a principled proposer's strategy is optimal.¹⁴ ■

The mini-ultimatum game, in its simplicity, shows many of the implications of our model regarding punishment. First, principled people punish -i.e., reject an offer- because they feel angry at violators of the E-norm. Second, angry responders trade off their desire for revenge and their material interest. Note that rejecting an offer costs $(1 - \rho)M$, that is, the amount of the offer. As ρ decreases, the cost of punishment increases, and that explains why principled responders do not reject very large *unfair* offers. The threshold depends crucially on the aggressiveness parameter α .

¹⁴This proposition holds for almost any ρ . More than one equilibrium exists if $\rho = \frac{1}{1+\alpha}$ ($\rho = 0$) because principled (selfish) responders are then indifferent between accepting or rejecting the unfair offer. The interested reader may easily find those equilibria.

The previous points are consistent with empirical evidence. Table 1 shows data reported in Slonim and Roth (1998) from one ultimatum game in which the stake size was 1500 Slovak Crowns (Sk), valued almost 48.5\$ at the exchange rate of that moment. For instance, 32.4% of all offers were in the offer range [40- 45) -i.e., each of them was larger or equal than 40% of the stake and smaller than 45% of the stake- and 4.9% of these offers were rejected. Consistent with our prediction, low offers are frequently rejected, and the probability of rejection tends to decrease as the offer increases. This result has been replicated in many other ultimatum game experiments.¹⁵

TABLE 1
SUMMARY OF SLONIM AND ROTH (1998)

Percentages of offers and rejections by range of offers

<u>Offer ranges</u>	<u>% Offers</u>	<u>% Rejections</u>
> 50%	7.2	0
= 50%	30.8	1.3
[45-50) %	6	0
[40-45) %	32.4	4.9
[35-40) %	5.2	0
[30-35) %	7.2	11.1
[25-30) %	3.2	37.5
<25 %	8	60

With respect to the proposer’s behavior, our model predicts that she will never offer more than half of the cake, which is basically consistent with experimental evidence. Moreover, the precise offer a proposer makes depends on her type, parameters α and γ , the size of the cake M , and priors μ . Let us consider each one separately.

The proposer’s type and parameter γ largely influence her degree of norm compliance. To make this point clear, assume for a moment that principled players may differ on the degree of internalization of the E-norm so that each one is characterized by a particular γ_i . In that case, principled proposers with a sufficiently large γ_i would choose equal sharing -note that proposition 6 still applies in this case. On the other hand, selfish and

¹⁵See Camerer [2003] or Güth [1995] for evidence on this. Note that one could easily introduce more heterogeneity regarding anger parameter α . If conveniently modelled, this idea could indeed explain why offers are scattered.

weak-willed principled proposers -i.e., those with a small γ_i - tend to choose meaner offers, if available.

Nevertheless, if the amount of money M at play is large enough, even a large parameter γ_i might not be enough to offset the material benefits of deviating from the fair sharing. In other words, the larger the size of the stake, the meaner (in *percentage*) the average offer. In any case, this should not be overemphasized: Experimental evidence shows that changes in stakes have small effect on proposals.¹⁶

Notice that, interestingly, our model predicts a positive correlation between the offers that *a same agent* would make in the dictator and the ultimatum games, specially if the stake is not big. Although we are not aware of any within-subjects experiment testing this, we can at least compare ultimatum and dictator game data coming from between-subjects designs. The two most important results are that offers are less concentrated in the dictator game than in the ultimatum game, and that average offer is smaller in the dictator game. Our model is consistent with those facts and explains them because of the impossibility to punish deviators in the dictator game.

Finally, and as we saw before, parameter α determines principled responders' acceptance threshold. The larger α is, the larger such threshold is. Well informed proposers who pretend to deviate from the E-norm should take this into account and, consequently, adjust their offers to their beliefs about α and μ . In fact, it is a robust experimental fact that there are almost no offers below $0.2M$. This seems to indicate that deviant proposers expect a high proportion of aggressive responders.

5.2 Intentions and Punishment: The Choice Set Matters

A principled player's utility function is *path-dependent* because its functional form depends on previous history. Hence, our model differs radically from consequentialistic models in which utility only depends on the material outcome of an interaction: One's feelings in two different games may differ *even* if the material outcome coincides.

Path-dependency is crucial to understand violence and conflict. One crucial idea in this respect is that *the choice set matters*: An action with equal material consequences may be perfectly right in one setting but not in another in which, due to a larger choice

¹⁶Notice incidentally that actual responders' thresholds change also modestly with stake changes. Evidence on those two points is surveyed in Camerer (2003, pp. 60-62).

set, the norms at work commend different behavior. To illustrate this, imagine one firm and a trade union setting wages: A wage increase that is *fair* during a recession may be completely insulting during a period of expansion -see Kahneman et al. (1986) for evidence on this. As a result, workers' reactions in each case -e.g., the probability of going to strike- may differ.

As another illustration, consider the mini-ultimatum *material games* represented at Figures 4a and 4b. Player 1 can either offer 'left' (l) or 'right' (r) and player 2 can accept (A) or reject (R) any offer. In both games, 'left' consists of an offer to give eight and two monetary units to player 1 and 2, respectively. In game (5/5), 'right' is an offer to share equally ten monetary units while in game (10/0) 'right' consists of a demand of the whole cake for player 1.

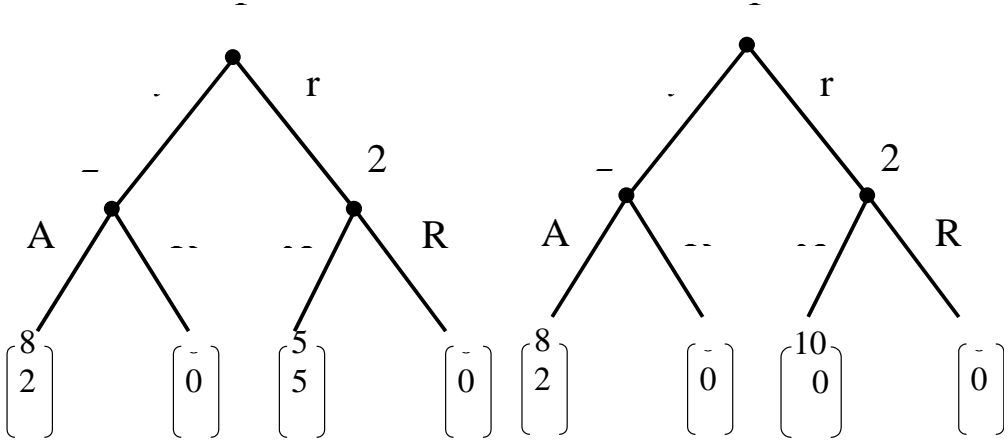


Figure 4a. (5/5) Material Game. Figure 4b. (10/0) Material Game.

Consider first player 2's behavior if he is offered 'left'. Falk, Fehr and Fischbacher (2003) find that 44.4% of the (8, 2) offers are rejected in game (5/5) while only 8.9% of those are rejected in game (10/0).¹⁷ Furthermore, proposers were able to anticipate the different rates of rejection. Around 30% of the proposers offer (8,2) in game (5/5) while almost 100% of the proposers offer (8, 2) in game (10/0).

In (5/5) game, our model predicts that principled responders reject left -assuming α is high enough- and accept right. The reason is simple: Offering (8, 2) in this game is an

¹⁷Brandts and Solà (2001) report similar results. See Fehr and Schmidt (2006) for more evidence and references on the topic.

unfair move because there exists a more fair offer -the equal sharing. Thus, offering (8, 2) activates anger and provokes rejection. Selfish responders, on the other hand, accept any offer. Therefore, a selfish *proposer*'s move depends on her initial expectation μ that the opponent is principled. She offers (8, 2) if μ is small enough and (5, 5) otherwise. Finally, principled proposers offer the equal split if γ is large enough.

Contrary to (5/5) game, the fairmax distribution of game (10/0) is (8, 2). Hence, offer (8, 2) is always accepted and this explains why offer (8, 2) is rejected at different rates in each game.¹⁸ On the other hand, an offer of (10, 0) is unfair and very cheap to punish so that it is always rejected by principled responders, whereas selfish ones are indifferent between accepting or rejecting it.

As another illustration of the influence of non-chosen alternatives, consider the mini-best-shot material game. Player 1 moves either 'left' (*l*) or 'right' (*r*). Player 2 observes her move and then either 'accepts' (*A*) or 'rejects' (*R*). Figure 5 shows its *material game* tree, in which $\rho M > (1 - \rho)M$ that is, $\rho > \frac{1}{2}$. A remarkable feature of this material game is that it has two fairmax paths: One leads to allocation $[\rho M, (1 - \rho)M]$, the other one to allocation $[(1 - \rho)M, \rho M]$. This largely drives our predictions.

¹⁸As we noted before, offer (8, 2) is rejected by a small minority of responders in game (10/0) so that our model is inconsistent with that result. May this phenomenon be due to inequity aversion? To analyze this point with a bit of detail, consider an *individual decision problem* in which the chooser must decide between two (self, other) material allocations: (2, 8) and (0, 0). Note that, from a consequentialistic point of view, this is exactly the same problem that a responder faces in the (10/0) game if the fair offer is made. However, and contrary to pure inequity aversion models, some experimental evidence shows that the rate of choice of allocation (0, 0) differs significantly in both cases -basically, *no* subject chooses (0, 0) in the non-strategic setting; see Charness and Rabin (2002) on this. All this seems to indicate that inequity aversion is not the force motivating rejection of the (8, 2) offer in game (10/0).

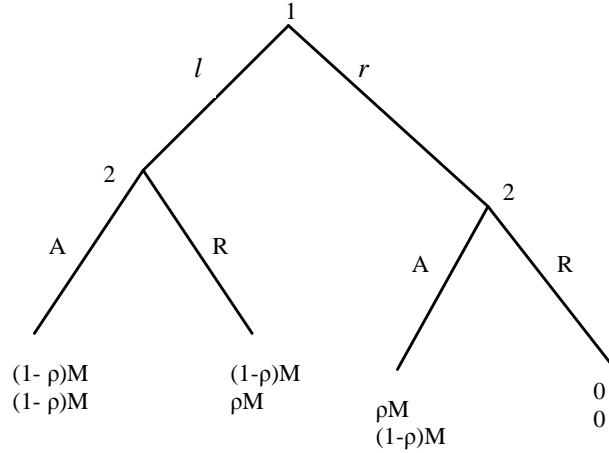


Figure 5. Mini-Best-Shot Material Game

Proposition 7 *The mini-best-shot game has a unique PBE strategy profile. Independently of their types, player 1 chooses ‘right’, and Player 2 rejects ‘left’ and accepts ‘right’.*

Proof. Player 2 accepts ‘right’ and rejects ‘left’ because that is consistent with the E-norm and maximizes material payoff. For the same reasons, player 1 offers ‘right’. ■

The mini-best-shot game shows how the model works in games with more than one fairmax path. The player that chooses at a decision node in which two or more fairmax paths diverge has a strategic advantage: She can choose the fairmax path that favours her most without making the opponents angry. Models of inequity aversion, on the contrary, predict that some of the responders will not accept ‘right’ because the ensuing distribution is disadvantageous for the responder.

Prasnikar and Roth (1992) study a best-shot game with a richer strategy space than the one we analyze here. Their results, however, are still consistent with our predictions. This is specially true concerning proposers. On the other hand, some responders prefer $(0, 0)$ than $[\rho M, (1 - \rho)M]$ -something that our model cannot explain. Does this indicate that they are inequity averse? For reasons that we mentioned before (consult the previous footnote), we doubt that. Furthermore, it is convenient to note that Prasnikar and Roth (1992) also study behavior in a comparable ultimatum game, and show that the rate of rejection of offer $[\rho M, (1 - \rho)M]$ is significantly *larger* in the ultimatum game.

The divergence in results cannot be explained by inequity aversion models because they assume consequentialistic preferences -i.e., the only thing that agents care about is how material resources are distributed, not how this distribution is achieved. Our

model, on the contrary, explains the divergence because offer $[\rho M, (1 - \rho)M]$ constitutes a deviation from the E-norm in the ultimatum game, where the equal sharing is feasible, but not in the best-shot game. Consequently, such offer makes the second mover angry in the ultimatum but not in the best-shot game.

To sum up, the whole set of alternatives is important because people determine what is fair by looking at that set. Consequently, an action may be fair in one context but not in another one in which a more fair move is feasible. Since anger is activated by unfair moves, it follows that punishment also depends on the initial set of alternatives.

5.3 Intentions and Punishment: Responsibility Matters

In our model, *responsibility* becomes an important variable because sanctions are directed only towards violators. Suppose, for instance, that vegetable production in a certain region has been minimal because of low irrigation, and think of two possible scenarios: In one, the cause of low irrigation was a heavy drought whereas in the other it was the incompetence of the agency in charge of the irrigation channels. Although distributional consequences -low agricultural incomes- may be identical in both cases, farmers are likely to anger at the agency in the latter scenario, thus generating conflict, but not in the former one. In general, unfair outcomes may be the result of Nature moves or third parties' choices. Economic crisis in little countries, to give another example, may be caused by policy choices made by big countries or international institutions. Citizens' response in this case is likely to be different that if the crisis is caused by bad economic policy at the domestic level.

It is possible to test in the lab whether responsibility matters or not. Suppose, for instance, that a computer generates randomly the proposer's offer in an ultimatum game. Since the proposer is not *responsible* of any deviation, the model predicts that no responder rejects (punishes) a randomly chosen offer x . Thus lower rejections rates are predicted in the computer treatment than in the typical, intentional treatment.

Blount (1995) was the first to provide experimental evidence on this regard.¹⁹ Our model is consistent with Blount's experimental data. Indeed, there is a significant and

¹⁹Blount's results are problematic because, among other reasons, subjects were deceived in one of the treatments. Further research shows, however, that Blount's qualitative results are not an artifact of the experimental design. See Fehr and Schmidt (2002) for a discussion.

substantial reduction in the acceptance thresholds of responders in the computer treatment.²⁰ Blount also studied rejection rates in case a third party chooses the proposer’s offer. Interestingly, acceptance thresholds in this condition did not differ significantly from those in the usual condition. Although this seems inconsistent with our model, we believe that it can be easily accommodated. When a third party chooses a sharing that favours either the proposer or the responder, that third party violates the norm of fairness. The responder may then ‘punish’ the third party by rejecting that unfair sharing -of course, rejection is more likely in case the split favours the proposer because then it is relatively cheaper.

The concept of responsibility is rather alien to consequentialistic models. Inequity aversion models as Fehr and Schmidt (1999) or Bolton and Ockenfels (2000), for instance, predict that a relatively well-off agent will always be sanctioned if punishment is cheap enough, *whatever* her previous behavior. As a result, they predict no change in rejection rates between the computer treatment and the typical, intentional treatment. On the contrary, models of intentions and type-based reciprocity predict some change.

5.4 Punishment is not a Means to Reduce Inequity

Models of altruistic motives like Charness and Rabin (2002) are unable to explain why people punish. On the contrary, inequity aversion models -Fehr and Schmidt (1999) and Bolton and Ockenfels (2000)- do provide a rationale: Punishment is a means to reduce disadvantageous *material* payoff inequality. For two-player games, this idea has a series of implications which we will contrast in what follows with ours.

First, inequity averse agents never punish an opponent that gets a *lower* payment than oneself. On the opposite, we predict that a principled agent will punish any transgressor (including disadvantaged ones) if punishment is cheap enough. As an illustration, consider the *material game* represented in Figure 6. Observe that, if given the choice, player 2 may punish the first mover by choosing R. Note also that player 2 gets a larger payoff than player 1 at any terminal node. Consequently, inequity aversion models predict that *any* second mover would play A if given the choice. A rational inequity averse first mover should then move *l* in order to get a larger material payoff and reduce disadvantageous

²⁰Nevertheless, there exists a very small proportion of actual responders that reject low offers in the random treatment. Again, see Fehr and Schmidt (2002) for a survey on this and related issues.

inequity.

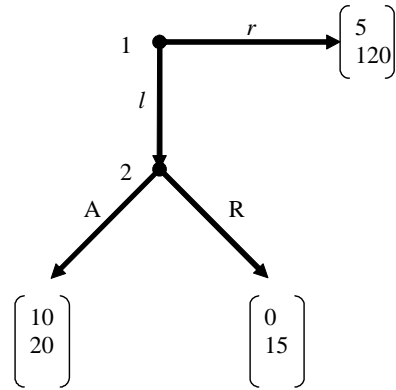


Figure 6. Punishing a disadvantaged opponent.

Unless the inequality parameter δ is very close to one, $(5, 120)$ is the unique fairmax distribution of this material game so that the E-norm prescribes to play r . Suppose then that player 1 violates the norm and plays l . If the second mover is principled and his anger parameter α is large enough, he will punish the first mover. That is, he will play R. As a result, a rational first mover may decide to play r if her priors μ are large enough.

Another prediction of inequity aversion models is that punishment never takes place if it is so costly that it does not reduce disadvantageous inequity. Suppose, for example, that reducing the opponent's payment in one monetary unit costs exactly one unit as well. Then no inequity averse agent would punish the other player. On the contrary, we predict that a very aggressive principled player -i.e., $\alpha = 1$ - would indeed punish a *transgressor*. Experimental evidence strongly supports our prediction -see Falk, Fehr and Fischbacher (2005) for details.

Finally, and to repeat one idea that was mentioned before, inequity aversion models predict some punishment towards an *advantaged* opponent if it is cheap enough, and independently of her previous behavior. Since our model predicts *no* punishment towards non-transgressors, independently of their relative status, it is clearly at odds with that idea. Experimental evidence seems to be at conflict too: In previously mentioned game Berk23 of Charness and Rabin (2002) subject B chooses between (B, other) allocations of dollars $(2, 8)$ and $(0, 0)$. Because punishment -i.e., choosing $(0, 0)$ - reduces inequity, inequity aversion models predict that a significant proportion of subjects choose $(0, 0)$. Contrary to that, *all* participants chose $(2, 8)$. Consult Fehr and Schmidt (2006) for more evidence.

5.5 Other Games

All the games we have analyzed so far are 2-player games. As an example of a multiple-player game, consider an ultimatum game with multiple proposers (the analysis here can be easily extended to an ultimatum game with multiple responders). That is, $n - 1 \geq 2$ sellers (proposers) make simultaneous price offers p_1, p_2, \dots, p_{n-1} to sell one unit of a good to a single buyer (responder) who demands only one unit of the good. The buyer can accept the offer she prefers or reject all of them. Assume that the responder values one unit of the good in V monetary units. Hence, the responder's monetary payoff if she accepts price offer p_i ($i \in \{1, 2, \dots, n - 1\}$) is $V - p_i$, whereas seller i 's income is p_i -unsuccessful sellers get zero money. Finally, all players get no money if the responder accepts no offer.

In order to study how players behave in equilibrium, assume that players' types are common knowledge (the analysis complicates if players' types are private information, and it does not add much insight). Note first that the E-norm selects just one fairmax path in this game: All sellers choose the same price ($V/2$) and the responder accepts one of these offers. This is clearly efficient and minimizes inequality between the worst-off player (an unsuccessful proposer) and the best-off one. Therefore, a principled responder gets a utility payoff $V/2$ if at least one proposer i makes a 'fair' price offer -i.e., $p_i = V/2$ - and the responder accepts it. On the contrary, she gets a payoff of $V - p_i - \alpha p_i$ if she accepts an unfair offer -i.e., $p_i \neq V/2$. In effect, the responder will be angry at i because he deviated from the E-norm, and moreover she will feel no remorse in that case because the E-norm allows any move once a deviation has been discovered. As a result, a principled responder's best response to any profile p_1, p_2, \dots, p_{n-1} is the following one: She accepts a fair offer if at least one proposer made such offer and if any unfair offer $p_i \neq V/2$ satisfies

$$p_i > \frac{V}{2(1 + \alpha)} = p^*,$$

and accepts the lowest unfair offer otherwise. On the other hand, it is clear that a selfish responder will always choose the lowest price.

Further, one can prove that at least one proposer will choose a zero offer in equilibrium if there are at least two selfish proposers (a selfish opponent would slightly undercut the lowest price offer otherwise; the argument here is similar to that in a Bertrand duopoly game). On the other hand, there exist multiple equilibrium *paths* if there is only one selfish proposer, or no selfish proposers. An equilibrium path in which all proposers make

a zero offer exists whatever the number of selfish proposers (to understand this, recall that principled types feel no remorse if everybody deviates from the norm). Another example appears when there is a single selfish proposer, *the responder is principled*, and $p^* - \varepsilon - \gamma(n - 2) < 0$ for any $\varepsilon > 0$. In this equilibrium path, the selfish proposer offers $p_i = p^*$ and all the principled proposers respect the E-norm (an analogous equilibrium exists if the responder is selfish and $V - \varepsilon - \gamma(n - 2) < 0$ for any ε arbitrarily close to 0, but the selfish proposer offers $p_i = V - \varepsilon$ in this case). To sum up, the outcome is largely undetermined if there are less than two selfish proposers, while the responder gets the whole surplus for sure otherwise. This latter result is largely consistent with the available experimental evidence, specially if subjects are allowed to play the game repeatedly; see the survey in Fehr and Schmidt (1999, p. 829).

In López-Pérez (2006) we have analyzed other games using a slight version of our model here, and compared predictions with existing experimental data. The model there explains, for instance, why people cooperate conditionally in the Prisoner's Dilemma, why first movers in a sequential social dilemma cooperate significantly more than players of a simultaneous dilemma, why cooperation in a Voluntary Contribution Mechanism (VCM) public good game depends on the expectation that sufficiently others will contribute as well, or why competitive markets induce principled people to behave as self-interested ones do.

The main difference between our model here and the model in López-Pérez (2006) is that principled players are assumed to care exclusively about the so-called Efficiency-cum-Maximin norm (or EM-norm), and not about the E-norm. The EM-norm is obtained just by defining a fairmax distribution as an allocation maximizing function $Q(x)$ -expression (3). As we have seen, this distinction is immaterial in two player games (which are our main focus in this paper), but it *could* generate different predictions in multiple player games. both norms select the *same* actions in a VCM public good game -that is, contributing the whole endowment to the public good- but they select different fairmax paths in competition games.

Consider again, for instance, the ultimatum game with multiple proposers. What does the EM-norm select to do in this game? As there are at least 2 sellers and only one of them can be succesful, it follows that, whatever players do, there will always be one seller who gets the minimum posible payment -i.e., zero. For that reason, the EM-norm selects

any price offer and commends always the responder to accept (because rejecting is not socially efficient). It is no wonder then that in equilibrium at least one player makes a zero offer, thus ensuring that the responder gets the whole surplus whatever the players' types (one can also show an analogous result for the ultimatum game with multiple *responders*).

6 Concluding remarks

We have shown that a large body of experimental evidence, including very different phenomena like generous and punishing behavior, may be explained by a relatively simple utility theory in which agents experience different emotional responses depending on how they and others act. Roughly speaking, our claim is that aggressive emotions like anger and moral emotions like guilt or shame are strong psychological forces that enforce reciprocity, understanding by that concept two things: (1) People adhere to norms if they expect others to respect them as well, and (2) people punish those who violate binding norms. Further, and since these emotions are activated by deviations from the norm, they induce path-dependent preferences.

Because it is very simple, the E-norm we propose is also too unrealistic. On one side, a norm that allows *any* move once a player transgresses it is a rather eccentric norm. Even if someone has committed a misdeed, actual norms of fairness still commend to be kind with those who previously respected them, and that may heavily restrict the set of decent choices. Another important issue is that societies have norms regulating revenge and punishment, something that the E-norm does not consider. For instance, proportionality concerns are widespread -i.e., many people believe that the punishment imposed on a deviator should be proportional to the damage that her deviation caused (*Lex Talionis*). Further, the E-norm has the problem that it is not strategic, that is, it makes prescriptions at any information set h without taking into account what the mover at h *expects* others are subsequently going to do -i.e., their intentions. Due to this, the E-norm may commend a move that, if the other player is selfish, ends up reaching a very unfair outcome.

To illustrate this, consider the *material game* tree represented at Figure 7. As $(3, 3)$ is the only fairmax allocation of this game, the E-norm commends player 1 to move l . Nevertheless, a selfish player 2 would then choose action A so reaching outcome $(0, 4)$,

which is most unfair. We believe that many people would argue that, unless one were sure enough that the opponent is a well-principled person who plays R, action r is at least as fair as action l . All this can be introduced in the model.

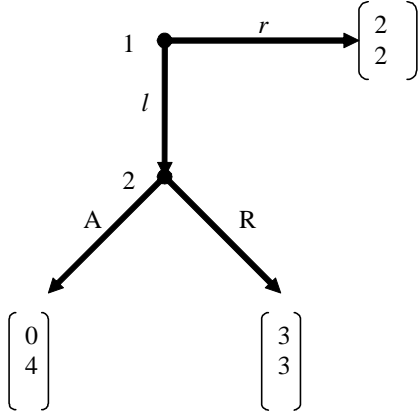


Figure 7. A Risky Move.

Assuming that the psychological cost triggered by a deviation does not depend on the particular deviation one makes is extreme. It seems more reasonable to assume that such cost grows with the *undeserved* harm our actions impose on the others. In particular, actions leading to unfair outcomes that favour the *opponent* should not induce any remorse at all -think of an ultimatum game proposer who offers the whole cake to the responder! Further, one could add a hypothesis of nonlinearity and some heterogeneity which would be useful, for instance, to replicate the fact that dictator game offers are usually scattered along the interval $[0, M/2]$.

This paper has concentrated on the study of fairness norms -i.e., norms regulating behavior in order to reach a fair distribution of material resources. Nevertheless, the definition of norm that we have given here is general enough to include many other types of norms. Think of norms regulating dressing, eating, or communication. For instance, parents instruct their children that telling lies is, most of the cases, a bad act that should embarrass them if performed. Accordingly, most of us feel badly when breaking that rule or anger at those who violate it. This emotional responses help to enforce sincere communication, and this can be easily accommodated within the setting offered by this paper. On the contrary, other models of social preferences and reciprocity are badly suited to explain these phenomena because they define a 'bad' action -if they define it at all- only by making reference to its expected material consequences; something that, obviously, communication does not affect.

To finish, the field of application of our model should not be restricted to the lab. Indeed, social norms and emotions strongly influence human action in many ‘real life’ settings like voting, law compliance, bargaining, team performance, and conflict, to cite a few. The next step should go in the direction of studying such influence -Lindbeck et al. (1999) is an example of this line of research.

7 Appendix:

As there are only two types of players in our model, it displays little heterogeneity. Although this is convenient for modelling reasons, it reduces sometimes the predictive power of the model. It is possible (and not very complex), however, to relax the assumption of common priors, and posit instead heterogeneous priors.²¹ We indicate in what follows one possible way to do it, and show how this could be used to better understand some experimental results.

Let μ_i denote the belief player i has about μ so that beliefs may be heterogeneous -i.e., $\mu_i \neq \mu_j$ for some $i \neq j$ - and mistaken -i.e., $\mu_i \neq \mu$. To simplify matters, we also assume that all players *believe* (maybe incorrectly) that priors are homogeneous and correct. More formally, if player i believes μ_i then i also holds the belief that μ_i is common knowledge. McKelvey and Palfrey (1992) called a hypothesis of this sort an *Egocentric model*.

This assumption is tractable and convenient because it does *not* require us to define a new solution concept. In effect, as player i believes that all players have common priors μ_i , we still predict that she will play according to a PBE of the game with common priors μ_i . To obtain behavioral predictions for any game, therefore, it suffices to find its PBEs as if priors were common.

As an application, consider how the Egocentric model affects our results in the Cournot duopoly game. We saw that the existence of the cooperative equilibrium required a sufficiently large common prior μ . When priors are heterogeneous, it is intuitive that the only principled players who will follow the norm are those who have a large enough

²¹Some experimental evidence supports this point of view. Palfrey and Rosenthal (1991), McKelvey and Palfrey (1992), Offerman et al (1996) and Sonnemans et al (1999) are some examples. Offerman et al (1996, pp. 838-839) remark that to explain their results "Not only is an altruism component needed in the utility function [...], but also an equilibrium concept which relaxes the assumption of accurate expectations".

prior. More precisely, if i is a principled firm with priors μ_i , she might produce q_F only if $4(3 - \mu_i)\sqrt{\mu_i G} \geq K$. Finally, heterogeneous priors also affect selfish firms' choices. To see that, observe that because θ depends positively on μ , a selfish agent i increases her choice $q_C = \theta q_F$ as her prior μ_i increases. In that way, we generate some behavioral heterogeneity from belief heterogeneity.

As another application, recall that we mentioned in the main text that participants in Cournot lab games usually converge towards the non-cooperative equilibrium when they play repeatedly. Although we do not pretend to offer any formal argument here (this would require a theory about how players update their beliefs), one could speculate that convergence is the result of mistaken priors: Participants who have initially large priors tend to update their beliefs and thus move to the non-cooperative equilibrium.

One might also apply the Egocentric model to our results in the Centipede Game. In the egocentric model, if player i has priors μ_i , then she plays according to a PBE of the game with common priors μ_i . As a consequence, a pessimistic selfish player -i.e., a player with a relatively low μ_i - will tend to terminate the game earlier than an optimistic one. This idea is completely consistent with the evidence and tests that McKelvey and Palfrey (1992) provide.

As a final example, consider the ultimatum game. If players have heterogeneous beliefs, our model predicts that selfish and principled proposers make *large* offers whenever her priors μ_i surpass a certain level, which is different for selfish and principled types. Now, in almost any experimental study the vast majority of offers is in the interval 40% – 50%. A possible interpretation of this phenomenon is that subjects come to the lab with rather large priors μ_i . In fact, the available evidence indicates that proposers tend to *overestimate* the actual proportion of people that reject unfair offers -i.e., parameter μ .²²

We finish with an esepculative remark on this last statement. From our knowledge, there is much experimental evidence (particularly from repeated games) that seems to indicate that average subjects overestimate parameter μ (is that prudence?). We believe that this issue would be an interesting topic for experimental research.

²²See Camerer (2003, p. 56) for a brief discussion of this point and some references. A natural question is whether more experienced subjects would adjust their strategies. The answer is positive: There is some evidence that offers slightly fall over time if repetition -with rematching- is allowed. Camerer (2003, pp. 59-60) also discusses this problem.

References

- [1] Andreoni, James and John H. Miller, 2002. "Giving according to GARP: An experimental test of the Consistency of Preferences for Altruism." *Econometrica*, 70(2), 737-53.
- [2] Blount, S., 1995. "When Social Outcomes Aren't Fair: The Effect of Causal Attributions on Preferences." *Organizational Behavior and Human Decision Process*, 63, 131-144.
- [3] Bolton, Gary E. and Axel Ockenfels, 2000. "ERC: A Theory of Equity, Reciprocity, and Competition." *American Economic Review*, 90(1), 166-93.
- [4] Bosman, Ronald and Frans van Winden, 2002. "Emotional Hazard in a Power-to-Take Experiment." *Economic Journal* 112, 147-69.
- [5] Brandts, Jordi and Carles Solà, 2001. "Reference Points and Negative Reciprocity in Simple Sequential Games." *Games and Economic Behavior*, 36, 138-157.
- [6] Burnham, Terence, 1999. "Testosterone and Negotiation: An Investigation into the Role of Biology in Economic Behavior." Harvard University, JFK School of Government.
- [7] Camerer, Colin F., 2003. *Behavioral Game Theory. Experiments in Strategic Interaction*. Russel Sage Foundation-Princeton University Press.
- [8] Charness, Gary and Brit Grosskopf, 2001. "Relative Payoffs and Happiness: An Experimental Study." *Journal of Economic Behavior and Organization*, 45, 301-328.
- [9] Charness, Gary and Matthew Rabin, 2002. "Understanding Social Preferences with Simple Tests." *Quarterly Journal of Economics*, 117, 817-869.
- [10] Cox, James C., and Cary A. Deck, 2005. "On the nature of reciprocal motives." *Economic Inquiry*, 43(3), 623-635.
- [11] Dufwenberg, Martin and Uri Gneezy, 2000. "Measuring Beliefs in an Experimental Lost Wallet Game." *Games and Economic Behavior*, 30, 163-182.

- [12] Dufwenberg, Martin and Georg Kirchsteiger, 2004. "A Theory of Sequential Reciprocity." *Games and Economic Behavior*, 47, 268-98.
- [13] Edgeworth, Francis Y., 1881. *Mathematical Psychics: An Essay on the Application of Mathematics to the Moral Sciences*. London: Kegan Paul.
- [14] Elster, Jon, 1999. *Alchemies of the Mind. Rationality and the Emotions*. Cambridge University Press.
- [15] Engelmann, Dirk; and Martin Strobel, 2004. "Inequality Aversion, Efficiency, and Maximin Preferences in Simple Distribution Experiments." *American Economic Review* 94(4), 857-869.
- [16] Falk, Armin; Ernst Fehr and Urs Fischbacher, 2005. "Driving Forces Behind Informal Sanctions." *Econometrica*, 7(6), 2017-30.
- [17] Falk, Armin; Ernst Fehr and Urs Fischbacher, 2003. "On the Nature of Fair Behavior." *Economic Inquiry*, 41(1), 20-26.
- [18] Falk, Armin and Urs Fischbacher, 2006. "A Theory of Reciprocity." *Games and Economic Behavior*, 54(2), 293-315.
- [19] Fehr, Ernst and Urs Fischbacher. "Third party Punishment and Social Norms." *Evolution and Human Behavior*, 25, 2004, 63-87.
- [20] Fehr, Ernst and Klaus Schmidt, 1999. "A Theory of Fairness, Competition and Cooperation." *Quarterly Journal of Economics*, 114(3), 817-68.
- [21] Fehr, E., and K. Schmidt, 2006. "The Economics of Fairness, Reciprocity and Altruism – Experimental Evidence and New Theories", in S. C. Kolm, and J. M. Ythier (Eds.), *Handbook of the Economics of Giving, Altruism and Reciprocity*, Vol. 1, Elsevier B. V.
- [22] Frijda, N., 1986. *The Emotions*. Cambridge University Press.
- [23] Geanakoplos, J., D. Pearce and E. Stacchetti, 1989. "Psychological Games and Sequential Rationality." *Games and Economic Behavior* 1, 60-79.

- [24] Gintis, Herbert, 2000. "Strong Reciprocity and Human Sociality." *Journal of Theoretical Biology*, 206, 169-179.
- [25] Güth, Werner, 1995. "On Ultimatum Bargaining Experiments—A Personal Review." *Journal of Economic Behavior and Organization*, 27, 329-344.
- [26] Holt, Charles A., 1995. "Industrial Organization: A Survey of Laboratory Research." in John H. Kagel and Alvin E. Roth, eds., *Handbook of Experimental Economics*. Princeton, NJ: Princeton University Press.
- [27] Kahneman, D., J. L. Knetsch, and R. Thaler, 1986. "Fairness as a Constraint on Profit Seeking: Entitlements in the Market." *American Economic Review*, 76, 728-741.
- [28] Kreps, D. M., P. Milgrom, J. Roberts, and R. Wilson, 1982. "Rational Cooperation in the Finitely Repeated Prisoner's Dilemma." *Journal of Economic Theory*, 27(2), pp. 245-52.
- [29] Levine, David K., 1998. "Modeling Altruism and Spitefulness in Experiments." *Review of Economic Dynamics*, 1, 593-622.
- [30] Lindbeck, A., S. Nyberg, J. W. Weibull, 1999. "Social Norms and Economic Incentives in the Welfare State." *The Quarterly Journal of Economics*, 114(1), 1-35.
- [31] López-Pérez, Raúl, 2005. "Guilt and Shame in Games." Unpublished paper.
- [32] López-Pérez, Raúl, 2006. "Introducing Social Norms in Game Theory." Working paper, University of Zurich.
- [33] McKelvey, Richard. D., and Thomas R. Palfrey, 1992. "An Experimental Study of the Centipede Game." *Econometrica*, 60, 803-836.
- [34] Offerman, Theo, 2002. "Hurting hurts more than helping helps." *European Economic Review* 46, 1423–1437.
- [35] Offerman, Theo, Joep Sonnemans, and Arthur Schram, 1996. "Value Orientations, Expectations, and Voluntary Contributions in Public Goods." *The Economic Journal* 106, 817-845.

- [36] Palfrey, T., and H. Rosenthal, 1991. "Testing Game Theoretic Models of Free-Riding: New Evidence on Probability Bias and Learning." in T. Palfrey, ed., *Laboratory Research in Political Economy*. Ann Arbor, MI: University of Michigan Press, 239-68.
- [37] Prasnikar, Vesna and Alvin E. Roth, 1992. "Considerations of Fairness and Strategy: Experimental Data from Sequential Games." *Quarterly Journal of Economics*, 107(3), 865-88.
- [38] Rabin, Matthew, 1993. "Incorporating Fairness into Game Theory and Economics." *American Economic Review* 83, 1281-1302.
- [39] Rosenthal, R., 1982. "Games of Perfect Information, Predatory Pricing, and the Chain Store Paradox." *Journal of Economic Theory*, 25, 92-100.
- [40] Sanfey, Alan G., J. K. Rilling, J. A. Aronson, L. E. Nystrom, and J. D. Cohen, 2003. "The Neural Basis of Economic Decision-Making in the Ultimatum Game." *Science*, 300, June 2003.
- [41] Slonim, Robert and Alvin E. Roth, 1998. "Learning in High Stakes Ultimatum Games: An Experiment in the Slovak Republic." *Econometrica*, 66, 3, 569-596.
- [42] Sonnemans, Joep, Arthur Schram and Theo Offerman, 1999. "Strategic Behavior in Public Good Games: When Partners Drift Apart." *Economics Letters* 62, 35-41.