# Intrinsic and Instrumental Reciprocity:

# An Experimental Study[*]

Luis Cabral

Erkut Y. Ozbay

*New York University, Stern*

*University of Maryland*

Andrew Schotter

*New York University, CESS*

April 17, 2012

Key Words: Reciprocity, Infintely Repeated Games, Veto Game

JEL Classification: C73, C92

## Abstract

In the context of an indefinitely repeated veto game, we devise an experiment to distinguish between alternative explanations of generous behavior (accepting negative payoffs): altruism, intrinsic backward-looking reciprocity, and instrumental forward-looking reciprocity. Our results are broadly consistent with the hypothesis that observed sacrifices are motivated by equilibrium selfish, forward-looking reciprocal behavior although we find a more subtle way in which past kindness affects behavior.

1

# 1    Introduction

Reciprocity is a significant part of the behavioral repertoire of humans (and other animals). People seem willing to sacrifice their material well being to help others. As summarized by Sobel (2005) such behavior comes in two basic varieties which he labels "intrinsic" and "instrumental" reciprocity. In intrinsic reciprocity, a kind (unkind) act by one social agent changes the preferences of the people he interacts with in such a way as to elicit kindness (unkindness) in response (see also Segal and Sobel (2007, 2008)). Intrinsic reciprocity is therefore preference based and likely to depend on the context of the game being played and the perceived intentions of the players.[1] In these theories, because reciprocity is motivated by a positive (negative) interpretation of the intentions of one's opponent, how one arrives at (or is expected to arrive at) a final payoff vector is an important component in determining whether behavior should be rewarded or punished. Such behavior (or its expectation) alters the weight that players put on the welfare of their opponents. In most intrinsic theories, see Rabin (1993) and Dufwenberg and Kirchsteiger (2004) Battigalli and Dufwenberg (2009) for example, when the game analyzed is not repeated, reciprocity results from the first and second order beliefs of the players about the intentions of the others which casts these models as psycholgical games. When games are repeated, as they are in this paper, it might make sense to think that subjects will look back at the previous play of their opponent in order to asses their kindnesses or perhaps their intentions and beleifs. This is, in fact, what we do here.

Other theories of reciprocity include altruism and the interdependent preference theories of Fehr and Schmidt (1999), and Bolton and Ockenfels (2000). These theories, differ from intrinsic models discussed above by ignoring the process through which final outcomes are determined and concentrating on the final distributions themselves. In other words, in these theories the preferences of agents are fixed and do not change in response to the behavior of others or one's perception of their intensions.

In contrast to intrinsic reciprocity, Sobel (2005) classifies reciprocity as instrumental if it is part of a repeated game strategy where agents sacrifice their short term gains in an effort to increase their long run (discounted) payoff. In such models, agents are capable of being perfectly selfish yet reciprocal behavior is observed as part of the equilibrium of the game. If Folk Theorems apply, a wide variety of behavior can emerge along with a wide variety of equilibrium outcomes all determined by selfish agents who are

---

[1] Many theories of reciprocty are cast as psycholgical games (see Rabin (1993), Dufwenberg and Kirchsteiger (2004), Battigalli and Dufwenberg (2009), Charness and Dufwenberg ( 2006), Celen, Blanco, and Schotter (2013)).

"forward looking" in the sense that they care about the impact of their actions today on the perceptions and actions of their opponent in the future. The logic of the Folk Theorem is the logic of instrumental reciprocity (see Rubinstein (1979), Fudenberg and Maskin (1986) and Abreu (1988), and more directly for our work here Cabral (2005)). [2].

In this paper we embed our experiment in an indefinitely repeated veto game of the type studied theoretically by Cabral (2005). In such veto games, in each of an infinite number of periods, Nature generates a pair of payoffs, one for each player. Although the sum of the players' payoffs is positive, one of the players may receive a negative payoff. Efficient equilibria thus require that players inter-temporally exchange favors, i.e., accept negative payoffs in some period with the expectation that such a favor will be reciprocated later in the interaction. An additional advantage of the repeated veto game is that, unlike most other repeated games, it admits a unique efficient equilibrium in the class of trigger strategy equilibria. We consider this equilibrium as the natural prediction of the selfish, rational behavior model and use its predictions as guide in our empirical section. We find significant support for the instrumental forward-looking explanation of reciprocity. [3]

---

[2] While indefinitely repeated games are a natural context within which to test theories of reciprocity, as Asheim and Dufwenberg (2003) point out, such reciprocity can be achieved even in finitely repeated Prisoners' Dilemma games. Hence it need not be a necessary condition. On a different point, Reuben and Seutens (2011) go even further and suggest that subjects may mistakenly apply rules of behavior best suited for long-term interactions outside the lab to tasks assigned them in an experiment that is only repeated a finite number of times.

[3] Our paper is not alone in suggesting that much of what looks like reciprocal or cooperative behavior can have instrumental motives. Reuben and Suetens (2011), using an indefinitely repeated prisoners' dilemma game, reach a conclusion similar to ours that a good deal of cooperative behavior can be explained strategically (see also Engle-Warnick and Ruffle (2006) and Engle-Warnick and Slonim (2006)). In a very clever design they have subjects play an indefinitely repeated prisoners' dilemma game using the strategy method where, just as in our paper, subjects are informed about when the last play of the game will occur. The game they look at is a dynamic game where player 1 moves first and then player 2 and both players write down a strategy of what they will do if the period they are in turns out to be the last period or not. The second player can also condition his action on whether the first player has cooperated or not. By looking at the strategies used by the players it is possible to identify their motives. They conclude that most cooperation observed is actually motivated by strategic considerations which are mostly reputation building by player 2.

Dreber, Fudenberg, and Rand (2011) also offer support that cooperative behavior in an infinitely repeated prisoners' dilemma game with noise is not motivated by inequality averse preferences but is rather payoff maximizing and competitive. In this game, subjects play an indefinitely repeated prisoners' dilemma game followed by a dictator game. They are also given a questionnaire after the experiment to elicit the motivation behind their behavior. The dictator game is run in order

The repeated veto game is of significant theoretical and applied interest. Cabral (2005) applies it to the problem of international merger policy, that is, the situation when a merger must be approved by multiple national authorities. A related context is that of interest rate setting by the European Central Bank, where individual member countries have veto power of changes on the interest rate level. An additional, closer to home, example is that of faculty recruitment, where different groups (e.g., micro and macro) have different preferences and hiring opportunities arise at an uneven rate.

All of these situations require that participants exchange favors over time. Hence, from the point of view of experimental economics, the indefinitely repeated veto game provides an excellent testing ground for the relative importance of altruism, intrinsic and instrumental reciprocity and selfishness as determinants of behavior. This is what we attempt to do in this paper.[4]

Methodologically, our paper makes several contributions since there are several features of our design that are new to the indefinitely repeated game literature. In particular, as mentioned above, it is one of the first papers to examine reciprocal behavior in indefinitely repeated games. Second, we present an innovation of some methodological use that ensures that no repeated interaction ends before at least some predetermined number of periods have transpired (in our experiment six) despite the fact that we use a probabilistic continuation rule to simulate discounting.[5] We do this by using a technique that makes the

to be able to correlate behavior in the prisoners' dilemma game with giving in the dictator game, a proxy for altruism.

Using the behavior of the subjects in the repeated prisoners' dilemma, their giving in the dictator game, and their answers to the questionnaire, Dreber et. al conclude that cooperation in repeated games is primarily motivated by long-term payoff maximization and that social preferences do not seem to be a major source of the observed diversity of play.

[4]While indefinitely repeated game settings are natural ones to use when testing for instrumental reciprocity, they are not necessary. A finite repeated game of the type examined by Kreps, Milgrom, Roberts and Wilson (1982) where the uncertainty about the existence of reciprocal types, may also lead to behavior that looks reciprocal but is actually instrumental. Further, Reuben and Suetens (2012) example of an experiment that identifies (rational) instrumental reciprocity and intrinsic reciprocity in a finite game context as is Muller, Sefton, Steinberg and Vesterlund (2008) where they examine strategic reciprocity by allowing subjects to use conditional strategies in a two period public goods game. In the experiment subjects play an indefinitely repeated game and can condition their strategy on whether the round of play is the last one or not. In the experimental game it is rational to use a forward-looking reciprocal strategy if the probability that the partner is intrinsically motivated is sufficiently high. Also, Muller et al. (2008) present experimental evidence for strategic reciprocity in a finitely repeated game. They let subjects play a 2-period public goods game, and ask subjects to submit choices in the second period, conditional on those in the first period.

[5]See Dal Bó and Fréchette (2011) for an excellent example of the approach where termination is stochastic. See also, Frechette and Yuksel (2013) for a comparison of the discounting method used here and several used by other investigators.

4

first six periods in any interaction deterministic with discounting yet allows these periods to blend into the stochastically ending portion of the experiment (periods 7 and above) in a behaviorally continuous manner. This allows us to make sure that we do not waste money on games that end "too soon". Third, two of our treatments have the added feature that when the last period is stochastically determined we inform the subjects that such period has arrived (see Reuben and Seutens (2011) for a similar treatment). In other words, while we use a stochastic stopping rule to end the indefinitely repeated game, in two of our four treatments we inform our subjects when the last period has arrived. In the context of our experiments, this allows us to identify whether their behavior up until that point was motivated by reciprocal or selfish motives.

In this paper we will proceed as follows. In Section 2 we will present the theory underlying indefinitely repeated veto games in the context of the experiment we conduct. In Section 3 we present our hypotheses while in Section 4, we present our experimental design. In Section 5 we present our results. Finally in Section 6 we offer the conclusions.

## 2    Theories of agent behavior

Our theoretical analysis is based on the following *repeated veto game*.[6] Two players interact over an infinite series of periods. Both players discount future payoffs according to the discount factor $\delta$. In each period $t$, Nature determines a *proposal*, a pair of payoff values $w_t = (w_{1t}, w_{2t})$ drawn from the set $S$ according to the c.d.f. $F(w)$, which we assume is smooth. Both players observe both values in $w_t$. Both players then simultaneously decide whether or not to approve the proposal $w_t$. If both players accept, then player $i$ receives payoff $w_{it}$. If at least one of the players rejects the proposal, then both players receive zero. Specifically, let $x_{it}$ be player $i$'s decision at time $t$, where $x_{it} = 1$ denotes approval and $x_{it} = 0$ denotes veto. Player $i$'s payoff in period $t$ is then given by

$$\pi_{it} = w_{it} \; x_{it} \; x_{jt}$$

Figure 1 illustrates a possible set $S$ (where for simplicity we drop the time component of the subscript

---

[6]See Cabral (2005) for a more extensive discussion of the repeated veto game and an application to international merger policy.

of $w$). All points in $S$ lead to a positive aggregate payoff.[7] We can consider three partitions of $S$. Points in region $A$ yield a positive payoff to both players. Points in region $D_i$ have the interesting property that (a) aggregate payoff is positive, (b) player $j$'s payoff is negative.
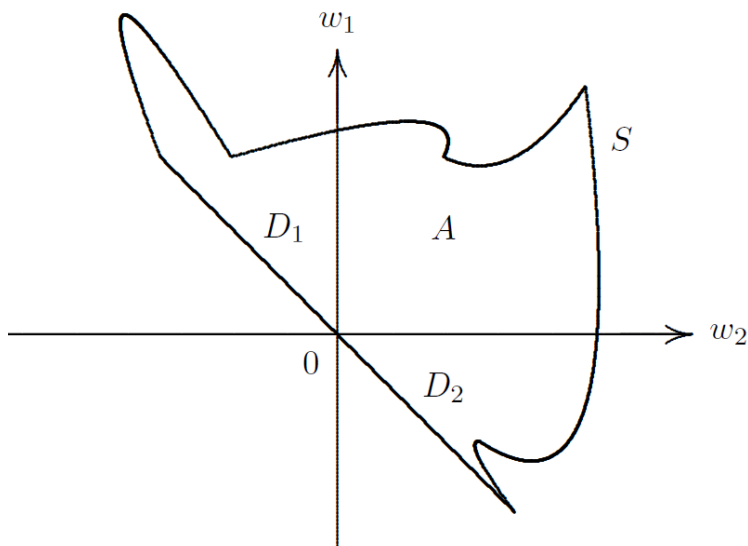


Figure 1: Payoff structure in a repeated veto game

It is straightforward to show that one equilibrium of this indefinitely repeated game would be to play a static Nash equilibrium in every period where each player rejects all negative payoffs for himself and accepts only positive payoffs no matter what offer is made to his opponent, or alternatively rejects all offers no matter whether they are positive or negative.[8] Experimental and anecdotal evidence suggest, however, that subjects are frequently "nice" to other players, that is, approve proposals yielding negative payoff for them but a positive aggregate payoff (that is, points in regions $D_i$). What theory can then explain the evidence? Our purpose in the present paper is to attempt to answer this question.

There are several reasons why outcomes do not correspond to the repeated play of static Nash equilibria. One first reason is that players care about other players' payoff: altruism or other regarding preferences. A second reason is that players follow some notion of reciprocity in their behavior: to the extent that

---

[7]Cabral (2005) considers the more general case when $S$ includes points with negative aggregate payoff.

[8]As we will discuss later, this second equilibrium is unlikely to be played especially since it is weakly dominated by the first. Still, we list it because it is a logical possibility.

their partner has been kind in the past, reciprocating such kindness yields positive utility. Finally, a natural explanation based on economic theory is that the outcome of cooperation corresponds to a Nash equilibrium of the repeated game which is different from the static Nash equilibrium; that is, given repetition, players might achieve an equilibrium whereby some points in regions $D_i$ get approved. We next develop each theoretical hypothesis in greater detail.

◊ Altruism and Other-Regarding Preferences. An explanation for "generous" behavior (proposals in region $D_i$ that are approved) is altruism, the idea that a player's utility includes the amount earned by the other player. This is captured by $\Phi(w_{it}, w_{jt}) : S \rightarrow \mathbb{R}$. Specifically, suppose that, in each period, each player's utility is given by his payoff plus a fixed positive coefficient $\alpha$ times the amount earned by the other player. Suppose, for the moment, that players are myopic, that is, they do not consider the continuation of the game. Such altruistic preferences imply the following definition.

**Definition 1 (altruism)** *Under myopic, altruistic play, $x_{it} = 1$ if and only if $\Phi(w_{it}, w_{jt}) > 0$, where* $\frac{\partial \Phi}{\partial w_{it}} > 0$ *and* $\frac{\partial \Phi}{\partial w_{jt}} > 0$.

Figure 2 illustrates the linear case, when $\Phi(w_{it}, w_{jt}) = w_{it} + \alpha\, w_{jt}$ (where $\alpha > 0$ is the coefficient of altruism). In this case, we expect all proposals to the Nash equilibrium above the $\ell_1$ and $\ell_2$ lines to be approved.
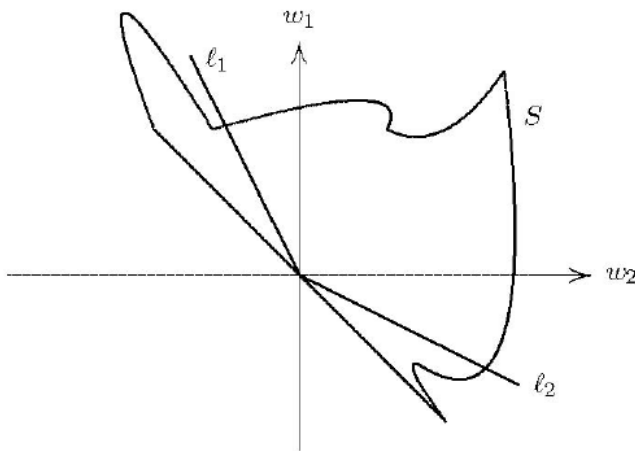


Figure 2: Altruistic, myopic equilibrium

Note that a similar result would hold if our subjects had various other types of other-regarding preferences such as those specified by Fehr and Schmidt (1999) and Bolton and Ockenfels (2000) since in both of these theories the decision to accept or reject an offer at any time t would depend both on one's own offer and that of one's opponent. It is important to note that if players consider the history of the game, and are, for example, inequality averse over payoffs accumulated over the game rather than payoffs from one period, this would lead to different predictions for inequality aversion than for altruism. Nevertheless, the players should take into account the other players payoff.

As we will see later, we need not restrict ourselves to myopic altruism since even in an indefinitely repeated game, if we assume that subjects use trigger strategies, the only efficient trigger-strategy equilibrium where people have non-selfish preferences involves subjects making their accept/reject decisions at each point in time on the basis of both offers and not just their own. This will not be the case when subjects have selfish preferences as will be true in the instrumental reciprocity model.[9]

◇ Intrinsic (backward-looking) reciprocity. An alternative explanation for "generous" behavior (proposals in region $D_i$ that are approved) is given by what we will call intrinsic reciprocity. Such explanations are backward looking since a player looks back at the previous behavior of his opponent, makes a judgement about how kind she has been, and then decides whether to accept a negative payoff based on how negative the payoff is and how kind the opponent has been.

The obvious question is how can we measure the kindness of a player? While there may be many ways to do this it is clear that whatever index one uses should take into account not only how much of a sacrifice (how negative a payment was accepted) a player has made in the past to help his opponent but also how much did a given sacrifice increase the opponent's payoff. For example, it is clear that player i is being kind to player j when he accepts a large negative amount. However, for any given sacrifice, we would consider player i as being more kind if the payoff of player j increased a lot rather than a little.

For that purpose, we define the kindness[10] of player i toward player j at time period $\tau$ as:

---

[9] In an experiment with a very different design than ours, Charness and Haruvy (2002) investigate whether they can separate altruistic, equity-based, and reciprocal motives in a labor market game. They find that reciprocity, distributive concerns, and altruistic considerations all play a significant role in players' decisions

[10] It is important to note that this kindness index is a generalization of Rabin's (1993) index by incorporating the amount of sacrifice and benefit in the index.

$$h_{i\tau}(x_{i\tau}|w_{i\tau}, w_{j\tau}) = [(x_{i\tau} - 1) - \frac{w_{i\tau}}{100}]\frac{w_{j\tau}}{100}I(w_{i\tau} < 0)$$

In this function $x_{i\tau}$ takes a value of 1 when an offer in period $\tau$ is accepted and zero otherwise while $I(w_{i\tau} < 0)$ is an indicator function taking a value of 1 when the offer to player i in period $\tau$ is negative (we are assuming that one does not exhibit kindness when one accepts a positive offer)[11]. Looking at the right hand side of our kindness index we notice two terms, one inside the square brackets and one outside. The term inside the brackets we will call the negative sacrifice component since it measures how much of a sacrifice player i is making when he accepts a given negative offer $w_{i\tau}$ in an effort to help player j. To understand this term, consider a given period $\tau$ and suppose that $w_{i\tau} = -60$. If player i accepts this proposal (so that $x_{i\tau} = 1$), then we say he is being kind to his partner to the tune of $.60 = (x_{i\tau} - 1)\text{-}\frac{w_{i\tau}}{100}$ where $x_{i\tau} = 1$ and $w_{i\tau} = -60$. The maximum value of kindness in a given period is therefore 1; it corresponds to the case when player i accepts a sacrifice of $-100$. Suppose however that the player rejects the same proposal of $-60$ (so $x_{i\tau} = 0$). We then say he is being kind (or rather, unkind) to the tune of $-.40 = (x_{i\tau} - 1) - \frac{w_{i\tau}}{100}$, where $x_{i\tau} = 0$ and $w_{i\tau} = -60$. Intuitively, the idea is that kindness corresponds to accepting large negative offers. In the limit when $w_{i\tau} = -100$ is accepted, we get one unit of kindness. Conversely, unkindness corresponds to rejecting offers that would imply a small sacrifice to player i. In the limit when $w_{i\tau} = 0$ is rejected, we get one negative unit of kindness (or one unit of unkindness). Accepting an offer that implies a small loss is not considered to be either kind or unkind. In the limit when $w_{i\tau} = 0$ is accepted, we get $(x_{i\tau} - 1) - \frac{w_{i\tau}}{100} = 0$. Likewise, rejecting an offer that would imply a large loss is not considered to be either kind or unkind. In the limit when $w_{i\tau} = -100$ is rejected we again get $(x_{i\tau} - 1) - \frac{w_{i\tau}}{100} = 0$.

To explain the second term, again suppose that $w_{i\tau} = -60$. If player i accepts this proposal (so that $x_{i\tau} = 1$) when $w_{j\tau} = 61$ or when $w_{j\tau} = 91$, we say he is being kind but the magnitude of his kindness will be higher when $w_{j\tau} = 91$ than when $w_{j\tau} = 61$ since his kind action will benefit player j more when $w_{j\tau} = 91$. Similarly, if he rejects this proposal, he is being unkind and again the magnitude will be higher when $w_{j\tau} = 91$ than when $w_{j\tau} = 61$.

---

[11]It is possible that we should consider positive offers since it may be that one way to exhibit kindness is to reject a positive offer as a way of preventing one's opponent, whom you care about, from trying to be kind to you by accepting a large negative offer. Such behavior is rare so we ignore it in our kindness index.

If players are reciprocal, we would expect a player's utility from approving a proposal to be increasing in his partner's past kindness. Hence in order to determine the kindness of player i toward player j up until period $\tau$, it might be natural to simply add up $h_{i\tau}(x_{i\tau}|w_{i\tau}, w_{j\tau})$ from periods 1 to $\tau - 1$. However, not all past periods are likely to be weighted equally in the mind of player j. He may give more recent periods an increased weight and place declining weights on the more distant past. To capture this fact we impose a set of declining weights on past actions of player i and formulate his cumulative kindness at period $\tau$ as follows:

$$k_{i\tau} = \sum_{t=1}^{\tau-1} \lambda^{(\tau-t-1)} h_{it}(x_{it}|w_{i\tau}, w_{j\tau}).$$

If players employ kindness to motivate their reciprocity then this leads to a different prediction regarding the outcome of the game:

**Definition 2 (intrinsic reciprocity)** *In an intrinsic reciprocity equilibrium, $x_{it} = 1$ if and only if $\Phi(w_{it}, k_{jt}, w_{jt}) > 0$, where $\frac{\partial \Phi}{\partial w_{it}} > 0$ and $\frac{\partial \Phi}{\partial k_{jt}} > 0$.*

In the particular linear case, a proposal is approved if and only if $w_{it} + \alpha\, k_{jt} + \beta\, w_{jt} > 0$, where $\alpha > 0$. In other words, if kindness matters for some $\lambda \in [0, 1]$ and for one of the indices, the coefficient of the kindness should be strictly positive.

◇ Equilibrium (forward-looking) reciprocity. Economists have understood for a long time that selfish, individual utility maximization is consistent with the observation of cooperative behavior when games are indefinitely repeated. While it is possible to define an infinite set of possible strategies in the repeated veto game (as in any repeated game), we concentrate, as is often the case, on trigger strategy equilibria. In fact, in the econometric analysis of our data we will try to identify whether our subjects employed the efficient equilibrium which, as we will demonstrate, can only be reached using trigger strategies. We do this not necessarily because we believe, a priori, that subjects will naturally gravitate to these types of strategies but rather to furnish a precise prediction from which we can evaluate behavior. If observed behavior differs qualitatively from the behavior consistent with efficient trigger strategies, then clearly we selected an incorrect benchmark for our data analysis. As we will see, however, the behavior of our subjects is broadly consistent with the use of trigger strategies while not precisely efficient ones. Further,

since there are an infinite number of Nash equilibria, if we did not select one for predictive purposes, then any behavior observed is likely to be rationalized by some Nash equilibrium, making the theory vacuous.

The idea of a trigger strategy equilibrium is to consider a "cooperative phase," where each player chooses $x_i^C(w_i, w_j)$; and a "punishment phase," where each player plays the static Nash equilibrium strategy $x_i^N(w_i, w_j)$; and the rule is that players choose $x^C(w_i, w_j)$ so long as all players have chosen $x^C(w_i, w_j)$ in previous periods.

Specifically, let $x_i^k(w_i, w_j) : S \rightarrow \{0, 1\}$ be an action mapping from the set of possible proposals into the set of possible actions in each period, where 1 corresponds to approval, 0 to veto; and $k = C, N$. With some abuse of notation, let $x_{it}$ be player $i$'s actual choice at time $t$. Define the following cooperation indicator:

$$c_t \equiv \begin{cases} 1 & if \ \ x_{i\tau} = x_i^C(w_{i\tau}, w_{j\tau}), \ \forall i, \tau < t \\ \\ 0 & otherwise \end{cases}$$

Then a trigger-strategy equilibrium is defined as follows.

**Definition 3** *A trigger-strategy equilibrium is characterized by strategies*

$$x_{it} = \begin{cases} x_i^C & if \ \ c_t = 1 \\ \\ x_i^N & if \ \ c_t = 0 \end{cases}$$

Notice that there is a Nash equilibrium strategy which is simply to approve a proposal if payoff is positive: $x_i^N(w_i, w_j) = 1$ iff $w_i \geq 0$. As we will see below, this is not the only Nash equilibrium that can be used in the punishment phase. However, depending on what Nash equilibrium is assumed to occur, we can sustain different payoffs in equilibrium. We are interested in characterizing those equilibria that are optimal given an out-of-equilibrium threat.

**Definition 4** *An optimal equilibrium is a trigger strategy equilibrium that maximizes the sum of the players' expected discounted payoffs.*

**Proposition 1** *(Equilibrium (Instrumental) Reciprocity) For a given threat to be used in the punishment phase, there exists a unique optimal equilibrium, and it is such that $x_i^C(w_i, w_j) = 1$ if and only if $w_i \geq -\ell_i$, where $\ell_i$ is increasing in $\delta$ and $\ell_i = 0$ if $\delta = 0$.*

A proof may be found in the appendix. Proposition 1 is illustrated by Figure 3.
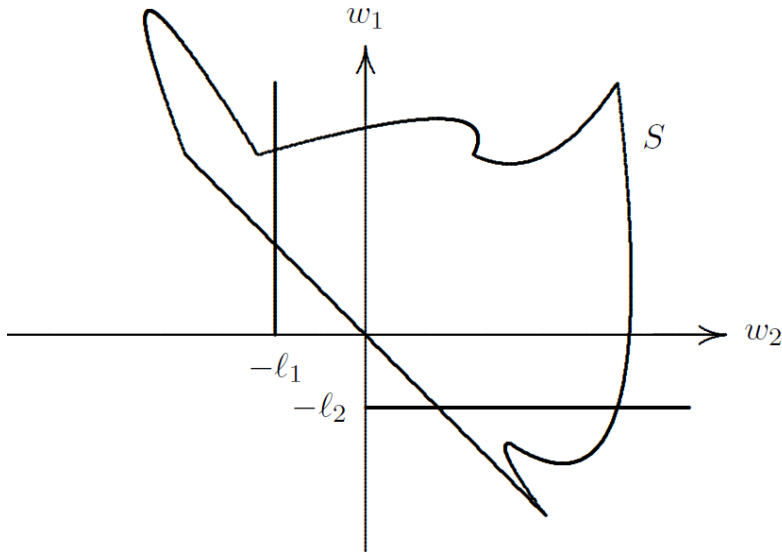
11

Figure 3: Optimal Threshold Equilibrium

Proposition 1 states that along the equilibrium path of the optimal equilibrium, all proposals in $S$ such that $w_1 > -\ell_1$ and $w_2 > -\ell_2$ are approved, and all the other ones are vetoed. Furthermore, for a given static Nash equilibrium to be used as a threat strategy in the punishment phase, there is only one pair $(\ell_1, \ell_2)$ that maximizes the sum of equilibrium payoffs.

Although in any Nash equilibrium of the one-shot version of this game, a player rejects any offer that gives negative payoffs to himself, there are multiplicity of Nash equilibria of this one-shot game. Some examples of the Nash equilibria of the one-shot game are when a player accepts a proposal if and only if his payoff is positive, or both players accept if and only if $w_1 > 10$ and $w_2 > 5$, or rejecting any offer.

During the punishment phase players may use any of these strategies. In the proof of the Proposition 1, we establish that whichever Nash strategy of the one-shot game is used as a threat, there exists a unique threshold strategy that maximizes the sum of the payoffs. Obviously, by using any of these strategies as a threat in the punishment phase, one supports several strategies as a part of the equilibrium of an indefinitely repeated game. For example, for the parameters used in the experiment, the threshold is -27 if accepting only positive offers is used as threat; it is -88 if rejecting any offer is used as a threat.[12]

---

[12]In the results section, we report thresholds as $-l$ rather than $l$ to emphasize that the subjects accept negative payoffs for themselves.

One may argue that despite the multiplicity of equilibria of the one-shot game, accepting only positive proposals is also the weakly dominant strategy and it might be unrealistic to think that the subjects will use the weakly dominated strategies as threats.

$\diamond$ Equilibrium (forward-looking) reciprocity with altruistic preferences. As we mentioned before we need not limit ourselves when discussing altruism to myopic behavior. The question then arises whether it is possible that when agents with such preferences interact over an infinite horizon they employ the same type of threshold trigger strategies as our selfish agents. The answer is no as long as we again restrict ourselves to efficient trigger strategies. In other words, if people are altruistic then in order to produce an efficient equilibrium in trigger strategies agents must take into account the payoffs of the agents they face no matter what threat is used. Since we find strong evidence that this is not the case, we again are presented with support for the notion that thresholds are used only by selfish agents.

These considerations yield the following proposition:

**Proposition 2** *(Equilibrium with Altruistic Preferences) If agents have altruistic preferences then no optimal trigger-strategy equilibrium exists in threshold strategies no matter what punishment threat is used and no matter what the functional form of the subjects' altruistic utility function is.*

A proof may be found in the appendix.

## 3   Hypotheses

The theory of instrumental reciprocity being tested here is characterized by two main features; thresholds and triggers. Thresholds characterize the cooperative phase while triggers characterize the punishment phase. If thresholds are employed by our subjects then we can rule out altruism or other-regarding preferences as a behavioral explanation since thresholds imply that the probability of accepting an offer in any round is  independent of the offer made to one's opponent, while altruism and other-regarding preferences suggest that the probability of accepting an offer depends on the offer of one's cohort. Hypotheses 1 and 2 concern these two features of our equilibrium.

**Hypothesis 1** *Thresholds: Subjects base their rejections of offers on the basis of a threshold above which offers are accepted and below which they are rejected. The probability that player i accepts a proposal is*

*increasing in player i's payoff and independent of j's payoff* .

The first part of this hypothesis obviously tests the threshold property of our model while the second part allows us to separate the impact of Instrumental Reciprocity from Altruism (or other-regarding preferences in general) since, as stated above, Instrumental Reciprocity with thresholds indicates that the rejection of an offer by subject i is independent of the offer made to subject j, while Altruism and other-regarding preference theories indicate that the probability of rejection depends on both offers. If we discover that including the consideration of an opponent's offer adds nothing to our ability to predict the probability that an offer is accepted, then we have provided evidence against altruistic and other-regarding preferences and in support of instrumental reciprocity.

Note that the fact that people use thresholds is only part of the demonstration that they were adhering to a forward looking reciprocal equilibrium since such an equilibrium also requires subjects to punish their opponent for the remainder of their interaction when they deviate. The punishment is to accept only non-negative offers. This yields the following hypothesis.

**Hypothesis 2** *Trigger Strategies: Subjects employ trigger strategies when playing the indefinitely repeated veto game.*

As Sobel (2005) has indicated, Intrinsic or preference-based reciprocity is a function of the previous behavior of one's opponent. If one's opponent has behaved in a kind manner, then such kindness changes the attitude of a decision maker towards his opponent by increasing the weight attached to his or her payoff in the decision maker's utility function. The opposite is true if the opponent behaves badly. Hypothesis 3 tests this Intrinsic Backward-Looking hypothesis and distinguishes it from both Altruism and Instrumental Reciprocity since neither of those theories are influenced by the past behavior of one's opponent. Instrumental reciprocity simply compares the current offer to the subject's threshold while Altruism looks at the value of both current offers. Neither looks at the previous behavior of one's opponent.

**Hypothesis 3** *Backward-Looking Reciprocity: The probability that player i accepts a proposal is increasing in player j's kindness index.*

While both Instrumental and Intrinsic Reciprocity exhibit reciprocal behavior, they do so for different reasons. With Intrinsic Reciprocity, a subject is rewarded for previous kindness while with Instrumental

14

reciprocity one cooperates (accepts a negative offer) in period t in the hope that such cooperation will be reciprocated in the future. This would imply that if it were announced to both players that their relationship would end in the current period, then we should not observe any subject accepting a negative offer in that period if he or she subscribed to the Instrumental or Forward Looking theory (since there is no future left), while a subscriber of the Intrinsic or Backward-Looking theory would reciprocate if the previous kindness level of his or her opponent were high enough. In other words, when there is no tomorrow there is no role for Forward-Looking reciprocity yet Backward-looking reciprocity may still operate.

**Hypothesis 4** *The probability that player i accepts a negative proposal in any period $t_i$ depends on whether the subject is informed that that period is the last period in the relationship he is in.*

Of these four hypotheses, Hypotheses 1-2 investigate Instrumental (Forward-Looking) Reciprocity. While Hypothesis 1 attempts to separate it from Myopic Altruism (and other behvioral theories that takes opponents payoff into consideration), Hypothesis 2 investigates whether trigger strategies were used. Hypotheses 3 and 4 try to identify whether Intrinsic (Backward Looking) or Instrumental (Forward Looking) behavior is what is observed in the data.

In the next two sections we describe the experiment we designed to test these various hypotheses (Section 4) and analyze statistically the data produced by the experiment (Section 5).

# 4   Experiment procedures and design

Our experimental design was created in an effort to test the theories described above. While we ran four treatments (to be described below) the experimental task engaged in by our subjects in each treatment was identical and can be described as follows. In each period, a pair of potential payoffs or offers $(w_1, w_2)$ is randomly determined. These values are uniformly drawn from the set determined by the following conditions:

$$-100 \leq w_i \leq 100, \ 0 \leq w_1 + w_2 \leq 100$$

This set is illustrated by the shaded area in Figure 4.

Figure 4: Experimental proposals generated.

Both players observe both values $(w_1, w_2)$. Players then simultaneously decide whether or not to approve the proposal. If both players approve the proposal, then each gets a payoff $w_i$. If at least one player vetoes the proposal, then both players receive 0.

The underlying model we test is one involving an indefinitely repeated game. Following the common practice, we implement the indefinitely repeated game as a repeated game that ends after each period with a continuation probability $\delta$ (hazard rate $(1 - \delta)$). In fact, for a risk-neutral player time discount and the probability a game will end are substitute elements in the discount factor.

This procedure creates an obvious practical problem, namely the possibility that the actual experiment lasts for a very short time (maybe just one period). In order to obviate this problem, we created a minimum time horizon, $T_{\min}$. Play of the game lasts at least $T_{\min}$ periods for sure; and for $t > T_{\min}$, we apply the hazard rate $1 - \delta$. Moreover, for $t < T_{\min}$ we introduce a payoff multiplier which decreases at rate $\delta$. This implies that, for a risk-neutral player, the future looks the same at every period of the game.

More generally, the formula for the multiplier $x_t$ is

$$
x_t = \begin{cases} \delta^{(t - T_{\min})} & if \ \ t \le T_{\min} \\[2mm] 1 & if \ \ t > T_{\min}, \end{cases}
$$

16

and the values used in the experiment are given in Table 1.

**Table 1:**

**Period Payoff Multipliers**

| Period | Multiplier |
|--------|------------|
| 1 | 3.05 |
| 2 | 2.44 |
| 3 | 1.95 |
| 4 | 1.56 |
| 5 | 1.25 |
| 6 | 1.00 |
| 7+ | 1.00 |

Note that in all periods before period 7, where stochastic discounting starts, the payoffs are multiplied by a constant greater than 1. For example, all payoffs earned in period 1 are multiplied by 3.05 making them more valuable than those earned in period 4, where the multiplier is only 1.56. The multiplier decreases until period 6 where it is equal to 1 and remains at that level from that point on. Note, however, that in period 7 the hazard rate $\delta$ takes over and it is in place from period 7 onward.

Table 2 presents the parameter values we used in our experiment. The minimum number of periods was set at $T_{min} = 6$ and the discount rate set at $\delta = .8$ (that is, after the sixth period the particular game ended with probability 20%). Each subject played this indefinitely repeated game ten times (that is, there were 10 rounds). Finally, the resulting equilibrium thresholds under the efficient equilibrium hypothesis is given by $-27$ (see the Appendix for the calculations).

In the experiment, 132 subjects were recruited from the undergraduate population at New York University via an electronic recruitment system that sends all subjects in the subject pool an e-mail offering them an opportunity to participate. Subjects played for francs which were converted into dollars at the rate of .6c per Franc.

Table 2: Experimental Parameters and Equilibrium Values

| Parameter | Value |
|---|---|
| Discount rate | 0.8 |
| Number of Rounds | 10 |
| Min number of periods ($T_{\min}$) | 6 |
| Equilibrium threshold | -27 |

## 4.1  Experimental Design

The experiment consisted of four treatments which differed by the matching protocol used and the level of information offered to the subjects in the last period of each round. In all treatments, subjects played ten rounds of an indefinitely repeated game. Subjects did not know ex-ante how many periods each round would last for, though they knew that there was a random continuation probability of $\delta = 0.8$. In two treatments (Treatments 2 and 4), subjects were randomly rematched with a new partner in each round, that is, after each indefinitely repeated game (randomly) ended, while in the other two treatments (Treatments 1 and 3) subjects stayed with their first round match for the entire 10 rounds of the experiment. Furthermore, in Treatments 2 and 3, before playing the last period of each round, subjects were told that the end-period had arrived, that is, that the period they were about to begin would be the last period of the current indefinitely repeated game. In the remaining two treatments, (Treatment 1 and 4), no such information was offered. In short, we conducted a 2 x 2 design with the treatments designated as FixedNotKnown (Treatment 1), RandomKnown (Treatment 2), FixedKnown (Treatment 3) and RandomNotKnown (Treatment 4) with 30, 28, 32 and 42 subjects, respectively.

We ran these treatments for two reasons. First, we used random matching because we feared that, with fixed matching, the ten rounds of the indefinitely repeated game might lose their independence. For example, subjects may build up a kindness reputation that spans across rounds. We do exploit the fixed matching protocol to demonstrate that while we fail to see strong evidence for intrinsic reciprocity within rounds of the experiment, across rounds we do find that subjects adapt the thresholds they use as a function of the kindness exhibited by the subjects they are repeatedly matched with. Second, we varied the last period information in order to compare the relative merits of the forward and backward reciprocity hypotheses.

# 5 Results

In this section we will present the results of our experiment. We will do this by testing each of the hypotheses stated above on the individual level using the data generated by our experiment. In the logit regressions, we controlled subject level fixed effects.

## 5.1 Hypothesis 1

To discuss Hypothesis 1 we will start with a descriptive analysis.

Figures 5a and 5b display the set of offers presented to two subjects in our Treatment 1, along with an indication of which offers were rejected dark (blue) diamonds and which were accepted light (purple) squares.

**Figures 5a and 5b: Individual Acceptance Behavior**



**Figure 5a (Subject 24)**



**Figure 5b (Subject 7)**

If Proposition 1 (and Hypothesis 1) is predictive of behavior, then in these graphs we should see a sharp division between offers that were accepted and those that were rejected with a rejection boundary separating the two that has an infinite slope. In other words, it should not be the case that the boundary between accepted and rejected offers has a negative finite slope.

As we can see, in Figure 5a this is certainly the case. For this subject (except for one observation) rejection behavior has the threshold property; offers above the threshold are accepted and those below are rejected regardless of the offer they imply for their opponent. Obviously, this was not  the case for all subjects, which is why we also present Figure 5b that shows the behavior of a subject whose attitudes appear to be more consistent with altruism since he seems willing to accept somewhat disadvantaged offers as long as they offer a large gain for his opponent. As our more formal regression analysis will indicate, these types of subjects are more the overwhelming exception than the rule.

Figures 6a and 6b (again from the Treatment 1) look at the data in another way. They present the acceptance behavior of subjects 19 and 13 in Treatment 1 over the 10 rounds of their participation in the game. On the horizontal axis we have the offer made to a given subject while on the vertical axis we measure two things. The first is a binary $\{0,1\}$ variable that takes a value of zero if an offer was rejected and a value of 1 if it was accepted. Second we measure the probability that a given offer is accepted using a logit regression where the binary accept/reject variable is regressed on a subject's own offer. If threshold behavior characterized a subject's behavior, then, when a simple logit function is fit to this data to explain acceptance behavior, our estimated logit regression should be a step function indicating that the probability of acceptance for offers below the step (threshold) is zero while it is one for offers above the threshold.

**Figures 6a and 6b: Acceptance Functions**



**Figure 6a**



**Figure 6b**

In Figure 6a we present our acceptance/rejection logit function for Subject 19 estimated by regressing his binary {0,1} response to his payoff offer. Note that Subject 19 beha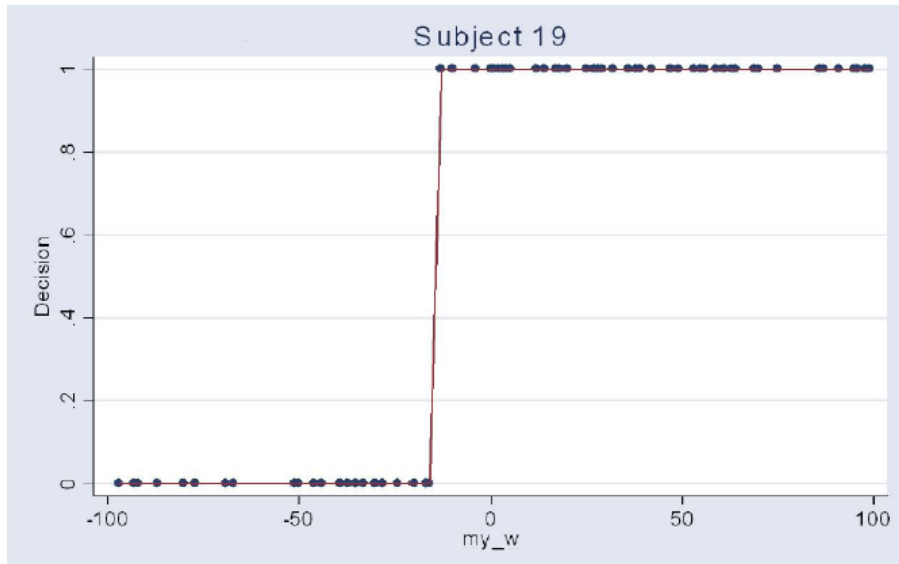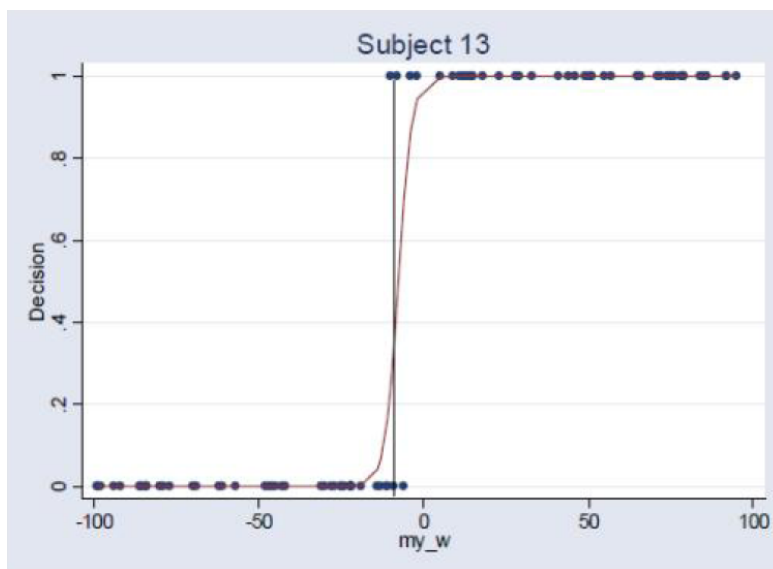ves exactly as a subject should if he or she was adhering to a strict threshold acceptance function. All offers below his threshold of -15

are accepted with probability 0 (rejected with probability 1) while those above the threshold are accepted with probability 1.

Subject 13 depicted in Figure 6b is a little different since his reject and accept regions for offers overlap. This means that this subject does not have a clear acceptance threshold. However, note that he is not far from perfect threshold behavior.

This discussion naturally leads us to look for a metric to use in assessing how far away a subject is from step function behavior and the pseudo threshold he is using. To do this we employ a very simple one which is to find the threshold which is such that we can fit a step-function to the data exactly by eliminating the minimal number of observations. To illustrate this, consider Figure 6b and Subject 13. From the logit acceptance function depicted there we see that, as opposed to Subject 19, this subject is not using a strict threshold acceptance function. This is true because the set of rejected and accepted offers overlap so there is not a clear separation between the sets of rejected and accepted offers. However, note that if we simply remove 2 observations from his data set (those to the right of the straight line on the bottom) we can establish a strict step function so this subject is 2 observations away from behaving as if he or she had a threshold strategy with a step at -10. Our metric then would award him a score of 2 and define his pseudo threshold as -10.

Tables 3a-3d in the Appendix presents, for each treatment, the estimated thresholds for each subject along with the number of observations that need to be eliminated to create perfect threshold behavior. This is followed by the percentage of the data not explained by these thresholds. Note that the exact threshold can not be uniquely defined by our procedure since there may be regions where no observations occur which straddle the actual threshold used. For that reason we provide two thresholds per subject (min and max) each of which can be used to estimate our threshold along with the mean threshold. In the remainder of the paper when we refer to a subject's threshold we will be referring to the mean stated in this table. [13]

As we can see from Tables 3a-3d, while not all subjects employed a perfect threshold strategy, many

---

[13] Another way of calculating the thresholds might be considering the logit regressions. Formally, the logistic function is $\exp(a + bx) / (1 + \exp(a + bx))$, so it takes the value of 1/2 when $a + bx = 0$. Since the threshold is the value of x for which subjects have 1/2 probability of taking either action, then the threshold $x^*$ can be found by setting $x^* = -a/b$, where a is the coefficient on the constant and b is the coefficient of the explanatory variable (own payoff). Tables 4a-4d presents these results. Our results are robust to these thresholds.

of them were in fact close to doing so in the sense that, on average, we only need to remove a few observations from each one in order to establish perfect threshold behavior. More precisely, note that over all rounds we only need to eliminate on average 5.43, 5.96, 5.46, and 5.5 observations from any subject in our four treatments respectively in order to establish perfect step-function behavior for him or her.[14] In addition, the maintained hypothesis that subjects used a threshold strategy is successful in explaining a large percentage of the data. For example, over all rounds the mean percentage of the data explained by our estimated thresholds are 94.38%, 94.05%, 94.32%, and 94.13% for treatments 1, 2, 3, and 4 respectively. This is strong support for the as if assumption that threshold behavior was operative.

These statistics actually under estimate how well threshold behavior fits our data. For example, from the logit regressions we will report later on in this section, our subjects naturally fall into two categories; those whose behavior can be explained exclusively with reference to one's own offer and those who take the offers of one's opponent into account as well.

Among the former group (constituting 102 of our 132 subjects) the mean number of observations that need to be eliminated in order to perfectly fit our rejection data with a step function is 3.2 while among those (25 subjects) who also care about one's opponent's offer (altruistic or intrinsically reciprocal types), the same number is 14.1. In other words, if we look only at the 77.3% of our subjects who exhibit strictly selfish behavior, our closeness index implies a closer fit.[15]

To test the second part of Hypothesis 1, we estimate a logit acceptance/rejection function for each subject $i$ by estimating the probability that $i$ accepts a offer $w_{it}$ given that $w_{jt}$ was offered to his pair member. We also include our previously defined opponent's kindness' variable, $k_{jt}$, in this regression, indicating the kindness of a subject's opponent up until the current period. In other words, we code the variable $a_{it}$ as a zero if the offer in period t was rejected and 1 if it was accepted, and we regress $a_{it}$ on $w_{it}, w_{jt}, k_{jt}$ and a constant. Since $k_{jt}$ is a function of $\lambda$, for each subject, we searched over $\lambda$, using values between 0 and 1 in steps of 0.10, to find that $\lambda$ which maximized the likelihood of the regression. These results are reported in Table 5a.

If Hypothesis 1 is accepted, then the coefficient on the $w_{jt}$ variable should be insignificantly different from zero while that of the $w_{it}$ should be positive and significantly so. Note that accepting Hypothesis 1

---

[14]The median of unexplained points are 3, 4, 3.5, and 2.5 in our four treatments respectively.

[15]Only 5 subjects can not be classifed at all.

is equivalent to rejecting the Myopic Altruistic or other-regarding preferences since those theories require a significantly positive coefficient on the $w_{jt}$ variable. Table 5a presents the summary of the results of our logit regressions run at the individual level for our four treatments.

**Table 5a and 5b here**

Additionally,Table 5b reports, for each treatment, the estimates of a pooled regression describing the probaility of accepting a negative offer given a proposal. As Table 5a and 5b clearly indicates, it appears that the probability of rejecting an offer for subjects is primarily a function of the offer they receive and not that received by their opponent. For example, over all subjects and all treatments of the 132 subjects who participated in our experiment,[16] 13 subjects had behavior that was perfectly described by thresholds in the sense that a step function (explaining rejections as a function of a subject's own offer) perfectly fit their data, 18 subjects had almost a perfect fit (only 1 unexplained point).For these subjects the estimated logit regression did not converge yet it is obvious that they only considered their own offer when contemplating rejections.[17] Including these subjects, 102 subjects had significant coefficients (at at least the 5% level) on their own offer variable, $w_{it}$ or had a perfect or almost a perfect fit. 25 also had significant coefficients on the $w_{jt}$ variable as well as $w_{it}$. None had a significant coefficient only on $w_{jt}$.In short, the primary determinant of rejection behavior seems to be one's own offer and not that of one's opponent. (We will discuss the coefficients on the kindness variable, $k_{jt}$, in a later section.)

These results present support for the threshold property of the Instrumental Reciprocity Hypothesis and for rejection of Myopic Altruism and other-regarding preferences since, under those hypotheses, a subject would have to take into effect his or her opponent's offer in determining the rejection and acceptance of an offer pair.

It is one thing to suggest that subjects behaved in a manner consistent with threshold strategies and yet another to suggest that they employed the theoretically optimal threshold of -27 in that strategy. Here our results suggest that while subjects did not use the theoretically optimal threshold in Random Matching treatments they did in the fixed matching treatments. More precisely, we calculated the weighted averages depending on how many observations had to be dropped.[18] Particularly, in Random Matching treatments

---

[16] If we restrict this Logit regression to only consider negative values for a subject's own offer, we get similar results.

[17] The regressions for these subjects are not, therefore, included in Table 5.

[18] For example, say Subject A's threshold is -20 and 1 out of 100 points need to be dropped; Subject B's threshold is -10

(Treatments 2, and 4) the weighted average of the thresholds were -11.90 and -5.93, respectively. By using a Wilcoxon signed rank test, we see that the weighted averages were significantly different than -27 (z=3.735, p = 0.0002; z=5.383, p=0.0000). It is important to note that, theoretically it is possible to show that those high thresholds can be explained by risk aversion. On the other hand, in the Fixed Matching treatments (Treatments 1 and 3), the weighted average of the thresholds were -17.97 and -19.23, respectively. For these thresholds we can not reject the hypothesis that they employed a threshold of -27 at 5% level by using a Wilcoxon signed rank test (z = 1.363, p=0.1728; z=1.758, 0.0787).

In conclusion, we have presented strong support for the idea that subjects employ threshold strategies. This result leads to rejection of the hypothesis that subjects were myopically altruistic or exhibiting other regarding preferences. However, we could not support the hypothesis that subjects employed the optimal thresholds across all treatments.

## 5.2   Hypothesis 2

The theory underlying these experiments relies on the use of trigger strategies with optimal thresholds. While we have offered support for the existence of threshold behavior, it is harder to detect whether our subjects used trigger strategies since punishments are only employed out of equilibrium. Given our data, however, it is hard to observe such out-of-equilibrium behavior. For example, one test as to whether triggers were employed would be to find a subject rejecting an offer that is better than what he/she had already accepted in an earlier period. This is true because if a threshold/trigger strategy is being employed, in the cooperative phase once an offer is accepted, all offers better than that one should be accepted as well. This would signal that the punishment phase had started. In our data, however, such occurrences are very rare (less than 1%) and, as a result, this test can not be used as evidence that triggers were employed.

Another feature of trigger strategies that should be observable in our data is the use of a common threshold for subjects who are paired together in the Fixed Matching Treatment. This is necessary since it must be commonly agreed upon as to when the punishment phase should be triggered. Hence, if optimal trigger strategies with the threshold property were used it would have to be the case that

---

and 5 out 100 points need to be dropped . Then, instead of taking the average of -20 and -10, we calculated [(-20)*99+(-10)*95]/(99+95)

our paired subjects used the same threshold during the experiment or at least converged to the same common threshold as time progressed. (Remember, for our experiment the optimal trigger is unique). The establishment of a common threshold takes time, however, at least for those subjects who do not have the ability to solve for the optimal equilibrium strategy. Hence, one explanation for the behavior of our subjects is that while they quickly learned to use a threshold strategy they had to interact over time to establish a common threshold upon which to base their trigger. If this is in fact the case, we should see the difference between the thresholds used by paired subjects in the Fixed matching treatments converge to zero over rounds. This is in fact what we see in Tables 6 (see Appendix Table 6a Table 6b for differences per pair),.

Table 6: Mean Difference in Pair Tresholds

| Treatment | All rounds | Rounds 1-5 | Rounds 6-10 |
|---|---|---|---|
| **FixedNotKnown** | 6.3 | 11.4 | 1.7 |
| **RandomKnown** | 11.9 | 11.8 | 12.2 |
| **FixedKnown** | 5.9 | 8.6 | 5.1 |
| **RandomNotKnown** | 14.5 | 15.8 | 14.2 |

Tables 6 present the average differences between the thresholds of paired subjects in FixedNotKnown and FixedKnown Treatments, and for RandomNotKnown and RandomKnown Treatments the average differences between one's threshold and the threhold of the group for the first and last five rounds of the experiment. As we can see, there is a general movement toward convergence in the thresholds used which is most pronounced in the FixedNotKnown treatment where, in rounds 6-10 the mean difference in the thresholds used was 1.7. This convergence lends support to the idea that our subjects were using trigger strategies but that it took time for our subjects to agree on a common threshold to serve as a trigger. Also as it can be seen in Table 6, converging as a group is harder than convergence as a pair, therefore the mean absolute difference from the mean is higher in the random matching treatments.

## 5.3 Hypothesis 3

If subjects subscribe to Intrinsic or Backward-looking reciprocity then the probability of accepting a negative offer in any period, t, should be positively related to the previous kindness of one's opponent up until period t-1. To test this hypothesis refer back to the regression reported in Table 5a where we

regressed our binary acceptance decision $a_{it}$ on a subject's offer in period t, $w_{it}$, his opponent's offer $w_{jt}$, and his opponent's kindness, $k_{jt}$ up to and including period t-1.

As is obvious from this table, we can strongly reject the hypothesis that subjects consider the previous kindness of their opponents when deciding whether or not to reject an offer. There were only 0, 6, 2 and 2 cases in which the kindness variable was significant at at least the 5% level in Treatments 1-4, respectively. Furthermore, if we look at the cases such that subject's own offer, partner's offer and kindness index are positive significant, there were only 2 cases in RandomNotKnown treatment, and zero cases in all other treatments.

The above results should not suggest that kindness reciprocity has no impact at all on behavior. We suspect that over time our subjects do respond indirectly to the kindness of their opponent by altering the threshold they use to accept and reject offers. To test this hypothesis we perform the following simple exercise. Using the data from our FixedNotKnown treatment, Treatment 1, we first divide the data into early (rounds 1-5) and late (rounds 6-10) rounds. We then correlate the change in thresholds used by our subjects from the first five to the last five rounds with the kindness of their opponents over the first five rounds. If our hypothesis is correct then we would expect a negative correlation between first-five-round kindness and the change in the thresholds used with more kindness observed in the first five rounds leading to lower (more negative) thresholds in the last five rounds. The correlation performed indicates that the relationship is negative, as it should be, with a correlation coefficient of -0.292 which is significant at the 5 % level. Hence, it would appear that kindness has an indirect impact of reciprocity - the kinder one's opponent is in the first five rounds the lower one's threshold is likely to be in the last five rounds. Such behavior may help to explain the convergence of thresholds noted on when discussing trigger strategies in Hypothesis 2.

## 5.4   Hypothesis 4

In our experimental design we run both fixed and random matching treatments with and without information. In the Known treatments we inform our subjects about the occurrence of the last round just before it is played. This allows a very natural test of whether subjects engage in backward (intrinsic) or forward looking (instrumental) reciprocity since, if subjects are backward looking, in the last round they should still be willing to reciprocate previous kindness with kind behavior by accepting a negative

offer while, forward looking behavior would rule out such a kind act since in the last period of a round subjects know they have no future together and hence the motivation to reciprocate is gone. Hence if Instrumental reciprocity were the guide to behavior we should see less negative offers being accepted in the last period of those treatments where information was full than in either the period just before the last or over all periods before the last. We expect to observe this behavior in the Random Matching treatments but not necessarily the Fixed Matching treatments since, in the Fixed Matching treatments, where people are rematched round after round, "last periods" lose their importance because subjects may still be willing to accept a low negative offer in a last period of round t in order to build a reputation that will be "reborn" in the round t+1 when they are rematched together. It is for this reason that we did the Random-Matching treatment in the first place.

As we see in Table 7a, our expectations were supported. Looking down column 1, we see that the fraction of negative offers accepted in the last period of the Random Matching Treatment was 0.112 while it was 0.191 for the period just before the last and 0.194 for all periods before the last ($p < 0.05$ in both cases). Note that, as expected, the same is not true for the Fixed Matching Treatment where the last period acceptance rates were 0.248 in the FixedKnown treatment and 0.218 in the FixedNotKnown treatment, and the difference is not significant. There are other comparisons which may be telling here as well. For example, we may want to compare the acceptance rates for subjects in the last periods of our two Random Matching treatments (Treatments 2 and 4) since both periods are last periods but in one that fact is known while in the other it is not. As we see, the acceptance rates are in fact lower with 0.112 of the offers being accepted when subjects know the offer was a last period offer while 0.159 were accepted when they did not know ($p < .05$).

|  |  | Random Known | Fixed Known | Random NotKnown | Fixed NotKnown |
|---|---|---|---|---|---|
| Last Period | Mean | 0.112 | 0.241 | 0.159 | 0.239 |
|  | SD | 0.317 | 0.430 | 0.367 | 0.428 |
|  | N | 98 | 116 | 157 | 113 |
| All Periods but Last | Mean | 0.194 | 0.225 | 0.154 | 0.197 |
|  | SD | 0.396 | 0.418 | 0.361 | 0.398 |
|  | N | 949 | 1033 | 1319 | 1026 |
| Next to Last Period | Mean | 0.191 | 0.248 | 0.141 | 0.218 |
|  | SD | 0.395 | 0.434 | 0.349 | 0.415 |
|  | N | 110 | 121 | 156 | 110 |

In order to control dependency of the aggregate data due to observations from same subjects, we ran logit regressions with subject fixed effects on the panel data where the left hand variable, "decision" was coded as a binary {0,1} variable where 1 denoted acceptance and 0 rejection. This variable was regressed on one of a set of dummy variables to be described below. We generated two dummy variables: *information* which assigns 1 if an observation comes from a treatment with Known (i.e. get information on whether the current period is the last period) and *lastperiod* which assigns 1 if an observation comes from the last period. By looking at the last period data only, in the Random Matching treatments (Treatments 2 and 4) we find that *information* has a significant effect on rejecting negative offers. This is not the case, however, if we look at the next to last period or all periods but the last one (see Table 7b). Additionally, in the RandomKnown treatment, *lastperiod* has a significant effect on rejecting negative offers ($coef. = -.078$, $SE = .038$, $N = 1047$, $p < .05$).

**Table 7b: Testing Negative Offer Acceptance,**

**Last and Not Last Rounds: Random Matching Treatments**

| | Random Matching All but Last | Random Matching Last Period | Random Matching Next to Last Period |
|---|---|---|---|
| Information | −0.05 | −1.00** | 0 |
| (SE) | (0.137) | (0.455) | (0.506) |
| Subject Fixed Effects | $YES$ | $YES$ | $YES$ |
| N | 2268 | 255 | 266 |
| $R^2$ | 0.17 | 0.37 | 0.30 |

$^{**}: p < .05$

One last comparison is interesting, and that is to compare the acceptance rates in the Random Matching Known and NotKnown treatments for all periods before the last. In other words, in these periods while the subjects in the Known treatment knew that that period's offer was not the last, subjects in the NotKnown treatment had to form a subjective estimate of the probability that that offer would be the last, an estimate that presumably increased as time went on and was positive in each period past the sixth. Under these circumstances we would expect that the acceptance rate in the Known treatment would be higher than in the NotKnown treatment since presumably subjects knew that these were still reputation building periods while subjects in the NotKnown treatment had a positive probability that this was the last period. Using data in the combined NotKnown and Known Random-Matching treatments for all periods but the last and regressing decision on *information (again* controlling subject fixed effects), supports the idea that acceptance rates are higher in the periods before the last when in the RandomKnown Treatment ($coef. = .25$, $SE = .087$, $N = 2268$, $p < .05$).

Our comments above lend support to the idea that most of the behavior we observed in this experiment, if it was reciprocal at all, was primary of the instrumental type. This is supported here by the fact that when subjects know they are in the last period of their interaction they tend to accept fewer negative offers while when they are not in the last period, but know that they will be informed when the last period comes, they accept more, presumably in an effort to keep their reputation alive.

## 5.5 Methodological contribution to indefinitely repeated games experiments

In laboratory experiments, indefinitely repeated games are induced by random termination. Using random termination may be costly, however, since some games may end quickly (even after only one period) and if they do they furnish little data for analysis. Because of this we introduce a novel method for our indefinitely repeated game experiments that allows collecting more data from each subject. To do this subjects first play the repeated game for fixed number, k, of periods (six in our experiment) with a discount factor and then play with random termination from period k+1 onward. The probability of termination is derived from the discount factor so that theoretically the two parts of the game "blend" into each other seamlessly. If this blending was, in fact, seamless, we should not observe any discrete change in the rejection probabilities of negative offers in the last (sixth) period of the deterministic phase and the first (seventh) period of the stochastic termination period. If we did, that would be evidence of a behavioral shift as we entered the stochastic phase of the round. To test this we pooled all of our data and compared the proportion of subjects accepting negative offers over two adjacent periods: the last period played with a deterministic discount factor (period 6) and the first period with a random termination (period 7) (conditional on that period not being the last in any treatment with Known information). What we find is that the fraction of negative offers rejected is practically identical across these two periods, 20% and 19.16% in the 6th and 7th periods, respectively, and these proportions are not significantly different (z=0.2706, p=0.7867). This result is what we hoped for since we wanted to smoothly bridge the transition between that portion of the game that was deterministic and that which was stochastic.

As a more formal approach to investigating whether acceptance behavior changes when we move across the boundary from periods 1-6 to periods 7 and beyond, we tested whether a structural break occurred in the estimated logit acceptance function between periods 1-6 and 7 and above, where the logit we were interested in had the {0,1} binary acceptance variable as a dependent variable and a subject's own offer (my_w) as the dependent variable using only those offers that were negative. To do this we first pooled all of our observations from all treatments. We then defined a dummy variable that takes a value of 0 if the observation came from period 1-6 and a value of 1 if it came from periods greater than 6. This dummy variable is entered as an independent variable and interacted with the intercept and slope coefficient in our logit estimation using a random effects specification for the error terms. This yields the following model (Model 1): $acceptance = \alpha + \beta_1(my\_w) + \beta_2 D + \beta_3 D(my\_w) + v_i + \epsilon_{it}$, where $\alpha$, $\beta_1$, $\beta_2$, and

$\beta_3$ are the coefficients to be estimated. We test the hypothesis that $\beta_2$ and $\beta_3$ independently are equal to zero as well as investigate whether they are jointly equal to zero. We do the latter by estimating the model with the restriction that $\beta_2 = \beta_3 = 0$ (Model 2) and performing a maximum likelihood ratio test. The results of this estimation are presented below.

| Table 8: Structural Break Regressions Results: Random Effects Logit | | | | |
|---|---|---|---|---|
| | Coef. | Std. Err. | z | P>\|z\| |
| my_w | .0648 | .0034 | 19.01 | 0.000 |
| Dummy | .2313 | 1696 | 1.36 | 0.173 |
| Dummy·my_w | .0009 | .0048 | 0.19 | 0.848 |
| Constant | .0658 | .1979 | 0.33 | 0.739 |
| N = 4811, Log likelihood = -1483.48, Prob > chi2 = 0.0000, | | | | |
| Log Likelihood Model 1 = -1483.48, | | | | |
| Log Likelihood model 2 = -1485.46, | | | | |
| chi2(2) =3.96, prob > 0.1373 | | | | |

As we can see, the result are consistent with the hypothesis that moving from a deterministic to a stochastic discounting regime after period 6 did not have any statistically significant impact of acceptance or rejection behavior. The $\beta_2$ and $\beta_3$ coefficients are both insignificantly different from zero indicating that there is no structural break in the acceptance function at period 6. In addition, the likelihood ratio test also indicates that $\beta_2$ and $\beta_3$ are jointly insignificantly different from zero.

In short, this regression lends support to the idea that our method of insuring a finite number of periods of play in our indefinitely repeated game did not alter the behavior of our subjects at the point where discounting became stochastic.

# 6    Conclusions

This paper has investigated the motives for reciprocal behavior in an indefinitely repeated veto game. In such games, in each of an infinite number of periods, Nature generates a pair of payoffs, one for each player. Although the sum of the players' payoffs is positive, with positive probability one of the players receives a negative payoff. In each period each pair member is asked to approve or reject the payoff pair.

If both subjects accept, then they receive the payoffs proposed, if one or more reject they both get zero. Clearly reciprocity in this game entails being willing to accept negative payoffs today with the hope that such generosity will be reciprocated in the future.

We consider this game to be a good vehicle to study reciprocity because the rationale for reciprocal behavior is obvious and the game is simple, despite the fact that it is indefinitely repeated. Following Cabral (2005) we designed an experiment whose purpose was to allow us to identify which one of two possible sources of reciprocity, intrinsic or instrumental, were most responsible for subject behavior.

Using some newly developed techniques to conduct indefinitely repeated games, our data supports the notion that in this indefinitely repeated game context, subject behavior is better described by theories of instrumental reciprocity but only to the extent that such reciprocity is part of a forward looking long run self-serving strategy. This is in distinction to intrinsic theories of reciprocity where reciprocal behavior is backward looking and exists to reward or punish previous kindness or unkindness. Despite this result, we find a number of ways that our subjects reciprocate kindness by sacrificing for opponents who have proven themselves to be kind in the past.

Finally, our results are consistent with the theory of veto games as presented in Cabral (2005) where optimal equilibrium behavior is characterized by a threshold for one's own payoff below which all offers are rejected but above which all offers are accepted regardless of the offer made to one's pair member.

# References

[1] ABREU, DILIP (1988), "On the Theory of Infinitely Repeated Games with Discounting," *Econometrica*, 56(2), 383–396.

[2] ASHEIM, GEIR AND MARTIN DUFWENBERG (2003). "Deductive Reasoning in Extensive Games," *Economic Journal,* 113, 305-325.

[3] PIERPAOLO BATTIGALLI AND MARTIN DUFWENBERG (2009). "Dynamic psychological games," *Journal of Economic Theory*, 144(1), 1–35.

[4] BOGACHAN CELEN, BLANCO, MARIANA, AND ANDREW SCHOTTER (2013), "On Blame-freeness and Reciprocity: An Experimental Study," mimeo.

[5] BOLTON, GARY, AND AXEL OCKENFELS (2000), "ERC: A Theory of Equity, Reciprocity, and Competition," *American Economic Review*, 90(1), 166–193.

[6] CABRAL, LUIS (2005), "An Equilibrium Approach to International Merger Policy," *International Journal of Industrial Organization,* 23, 739–751.

[7] CHARNESS, GARY AND HARUVY, ERNAN (2002) "Altruism, equity, and reciprocity in a gift-exchange experiment: an encompassing approach", *Games and Economic Behavior,* 40 203–231

[8] CHARNESS, GARY, AND MARTIN DUFWENBERG (2006), "Promises & Partnership," *Econometrica*, 74, 1579-1601.

[9] DAL BO, PEDRO AND GUILLAUME R. FRECHETTE (2011) "The Evolution of Cooperation in Infinitely Repeated Games: Experimental Evidence," *American Economic Review*, 101.

[10] DREBER, ANNA, FUDENBERG DREW AND RAND, DAVID (2011) "Who Cooperates in Repeated Games: The Role of Altruism, Inequity Aversion, and Demographics", mimeo, Yale University.

[11] DUFWENBERG, MARTIN, AND GEORG KIRCHSTEIGER (2004), "A Theory of Sequential Reciprocity," *Games and Economic Behavior*, 47, 268–298.

[12] ENGLE-WARNICK, JAMES, AND BRADLEY RUFFLE (2006), "The Strategies Behind Their Actions: A Method To Infer Repeated-Game Strategies And An Application To Buyer Behavior," Departmental Working Papers, 2005-04, McGill University, Department of Economics.

[13] ENGLE-WARNICK, JAMES, AND ROBERT SLONIM (2006), "Learning to Trust in Indefinitely Repeated Games," *Games and Economic Behavior*, 54, 95-114.

[14] FEHR, ERNST AND KLAUS SCHMIDT (1999), "A Theory of Fairness, Competition, and Cooperation," *Quarterly Journal of Economics*, 114(3), 817–868.

[15] FRECHETTE, GUILLAUME R. AND SEVGI YUKSEL, (2013) Infinitely Repeated Games in the Laboratory: Four Perspectives on Discounting and Random Termination.

[16] FUDENBERG, DREW, AND ERIC MASKIN (1986), "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information," *Econometrica*, 54(3), 533–554.

[17] KREPS, DAVID, PAUL MILGROM, JOHN ROBERTS, AND ROBERT WILSON (1982). "Rational Cooperation in the Finitely-Repeated Prisoners' Dilemma". Journal of Economic Theory 27 (2): 245-252.

[18] MULLER, L.,SEFTON, M., STEINBERG, R.AND VESTERLUND, L (2008) - Journal of Economic Behavior and Organization , vol. 67, n° 3-4, pp. 782-793.

[19] RABIN, MATTHEW (1993), "Incorporating Fairness into Game Theory and Economics," *American Economic Review*, 83, 1281–1302.

[20] RUBINSTEIN, ARIEL (1979), "Equilibrium in Supergames with the Overtaking Criterion," *Journal of Economic Theory*, 21, 1–9.

[21] REUBEN, ERNESTO, AND SIGRID SEUTENS (2012). "Revisiting Strategic versus Non-Strategic Cooperation," *Experimental Economics,* 15, 24-43.

[22] REUBEN, ERNESTO, AND SIGRID SEUTENS (2011), "Maladaptive Reciprocal Altruism," mimeo, Tilburg University.

[23] SEGAL, UZI AND JOEL SOBEL (2007) "Tit for Tat: Foundations of Preferences for Reciprocity in Strategic Settings," *Journal of Economic Theory*, 136, 197–216.

[24] Segal, Uzi and Joel Sobel (2008), "A Characterization of Intrinsic Reciprocity," *International Journal of Game Theory,* 36, 571–585.

[25] Sobel, Joel (2005), "Interdependent Preferences and Reciprocity," *Journal of Economic Literature,* 43, 392–436.

**7**

# Appendix A

**Tables 3a-3d: Thresholds**

Table 3a: Thresholds - Treatment 1 FixedNotKnown

| Subject | Max Threshold | Min Threshold | Mean Threshold | # of points unexplained | % of points unexplained |
|---|---|---|---|---|---|
| 1 | -9 | -11 | -10 | 1 | 1.0 |
| 2 | -9 | -9 | -9 | 2 | 1.8 |
| 3 | -83 | -100 | -91.5 | 6 | 5.8 |
| 4 | -1 | -8 | -4.5 | 3 | 3.0 |
| 5 | -6 | -6 | -6 | 2 | 2.3 |
| 6 | -15 | -29 | -22 | 37 | 38.1 |
| 7 | -27 | -27 | -27 | 6 | 5.8 |
| 8 | -8 | -8 | -8 | 4 | 4.7 |
| 9 | 11 | -21 | -5 | 15 | 19.2 |
| 10 | -42 | -43 | -42.5 | 7 | 5.3 |
| 11 | 2 | -2 | 0 | 2 | 1.5 |
| 12 | -35 | -42 | -38.5 | 3 | 3.1 |
| 13 | -4 | -10 | -7 | 2 | 2.3 |
| 14 | -19 | -23 | -21 | 4 | 4.1 |
| 15 | -4 | -4 | -4 | 5 | 6.4 |
| 16 | 0 | -1 | -0.5 | 5 | 3.7 |
| 17 | -43 | -45 | -44 | 13 | 9.9 |
| 18 | -18 | -26 | -22 | 5 | 4.6 |
| 19 | -13 | -15 | -14 | 0 | 0.0 |
| 20 | -31 | -38 | -34.5 | 6 | 6.2 |
| 21 | 0 | -1 | -0.5 | 2 | 1.6 |
| 22 | -97 | -100 | -98.5 | 16 | 19.5 |
| 23 | 4 | 3 | 3.5 | 2 | 2.2 |
| 24 | -4 | -12 | -8 | 1 | 0.9 |
| 25 | -17 | -20 | -18.5 | 3 | 3.1 |
| 26 | 2 | -2 | 0 | 2 | 2.2 |
| 27 | 1 | -1 | 0 | 0 | 0.0 |
| 28 | -7 | -17 | -12 | 1 | 0.9 |
| 29 | 5 | 1 | 3 | 2 | 2.1 |
| 30 | -18 | -29 | -23.5 | 6 | 7.3 |
| | | | | | |
| Mean | -16.17 | -21.53 | -18.85 | 5.43 | 5.62 |
| Median | -8.5 | -13.5 | -9.5 | 3.0 | 3.1 |

Table 3b: Thresholds - Treatment 2 RandomKnown

| Subject | Max Threshold | Min Threshold | Mean Threshold | # of points unexplained | % of points unexplained |
|---|---|---|---|---|---|
| 31 | -16 | -16 | -16 | 10 | 8.7 |
| 32 | -19 | -41 | -30 | 4 | 5.1 |
| 33 | -30 | -40 | -35 | 6 | 6.6 |
| 34 | 2 | -6 | -2 | 10 | 9.3 |
| 35 | -6 | -26 | -16 | 3 | 3.6 |
| 36 | -14 | -15 | -14.5 | 4 | 3.6 |
| 37 | -1 | -3 | -2 | 13 | 14.1 |
| 38 | 2 | 0 | 1 | 1 | 0.9 |
| 39 | -70 | -76 | -73 | 18 | 15.5 |
| 40 | 0 | -1 | -0.5 | 1 | 0.9 |
| 41 | -16 | -21 | -18.5 | 5 | 5.4 |
| 42 | 2 | 1 | 1.5 | 0 | 0.0 |
| 43 | 1 | 1 | 1 | 4 | 3.8 |
| 44 | 2 | 0 | 1 | 0 | 0.0 |
| 45 | -7 | -8 | -7.5 | 2 | 2.1 |
| 46 | 2 | -2 | 0 | 4 | 3.9 |
| 47 | -14 | -22 | -18 | 14 | 12.1 |
| 48 | 3 | 1 | 2 | 2 | 2.0 |
| 49 | -2 | -2 | -2 | 2 | 1.7 |
| 50 | -7 | -8 | -7.5 | 9 | 8.7 |
| 51 | -12 | -14 | -13 | 10 | 9.6 |
| 52 | -7 | -7 | -7 | 4 | 4.3 |
| 53 | -8 | -40 | -24 | 7 | 7.9 |
| 54 | -33 | -51 | -42 | 11 | 10.8 |
| 55 | -10 | -10 | -10 | 6 | 6.7 |
| 56 | -3 | -31 | -17 | 15 | 17.0 |
| 57 | 9 | 6 | 7.5 | 1 | 1.1 |
| 58 | -8 | -10 | -9 | 1 | 1.2 |
|  |  |  |  |  |  |
| Mean | -9.29 | -15.75 | -12.52 | 5.96 | 5.95 |
| Median | -7.00 | -9.00 | -8.25 | 4.00 | 4.69 |

Table 3c: Thresholds - Treatment 3 FixedKnown

| Subject | Max Threshold | Min Threshold | Mean Threshold | # of points unexplained | % of points unexplained |
|---|---|---|---|---|---|
| 59 | 6 | -4 | 1 | 6 | 5.5 |
| 60 | -37 | -58 | -47.5 | 8 | 9.3 |
| 61 | -6 | -34 | -20 | 20 | 19.6 |
| 62 | -6 | -19 | -12.5 | 21 | 23.1 |
| 63 | -31 | -31 | -31 | 8 | 9.3 |
| 64 | -4 | -10 | -7 | 3 | 3.3 |
| 65 | -19 | -26 | -22.5 | 1 | 1.0 |
| 66 | -95 | -96 | -95.5 | 3 | 3.2 |
| 67 | 0 | -5 | -2.5 | 0 | 0.0 |
| 68 | -10 | -17 | -13.5 | 2 | 2.0 |
| 69 | -5 | -13 | -9 | 9 | 11.7 |
| 70 | 2 | -1 | 0.5 | 3 | 2.8 |
| 71 | 3 | -5 | -1 | 2 | 2.0 |
| 72 | 0 | 0 | 0 | 1 | 1.1 |
| 73 | 4 | -2 | 1 | 0 | 0.0 |
| 74 | -10 | -14 | -12 | 4 | 4.1 |
| 75 | -99 | -100 | -99.5 | 2 | 2.1 |
| 76 | -2 | -19 | -10.5 | 4 | 3.9 |
| 77 | -54 | -57 | -55.5 | 12 | 11.3 |
| 78 | -18 | -18 | -18 | 6 | 7.8 |
| 79 | -7 | -7 | -7 | 2 | 2.2 |
| 80 | 0 | 0 | 0 | 7 | 6.4 |
| 81 | -13 | -13 | -13 | 21 | 20.2 |
| 82 | -38 | -53 | -45.5 | 6 | 5.0 |
| 83 | 4 | -2 | 1 | 2 | 1.9 |
| 84 | -83 | -83 | -83 | 4 | 3.3 |
| 85 | -2 | -5 | -3.5 | 0 | 0.0 |
| 86 | 5 | 5 | 5 | 2 | 2.2 |
| 87 | 0 | -5 | -2.5 | 6 | 7.1 |
| 88 | -3 | -4 | -3.5 | 3 | 2.9 |
| 89 | -1 | -1 | -1 | 1 | 1.0 |
| 90 | 6 | 1 | 3.5 | 6 | 6.6 |
|  |  |  |  |  |  |
| Mean | -16.03 | -21.75 | -18.89 | 5.47 | 5.68 |
| Median | -4.50 | -11.50 | -8.00 | 3.50 | 3.32 |

Table 3d: Thresholds - Treatment 4 RandomNotKnown

| Subject | Max Threshold | Min Threshold | Mean Threshold | # of points unexplained | % of points unexplained |
|---|---|---|---|---|---|
| 91 | -6 | -17 | -11.5 | 6 | 7.4 |
| 92 | 3 | -4 | -0.5 | 4 | 4.1 |
| 93 | -13 | -22 | -17.5 | 2 | 2.4 |
| 94 | -19 | -35 | -27 | 3 | 3.1 |
| 95 | 0 | -1 | -0.5 | 1 | 1.1 |
| 96 | 5 | -1 | 2 | 1 | 1.4 |
| 97 | 4 | -1 | 1.5 | 1 | 1.2 |
| 98 | 5 | 4 | 4.5 | 3 | 2.8 |
| 99 | 100 | 100 | 100 | 25 | 26.9 |
| 100 | -18 | -35 | -26.5 | 1 | 1.1 |
| 101 | -1 | -4 | -2.5 | 6 | 5.5 |
| 102 | 1 | -2 | -0.5 | 0 | 0.0 |
| 103 | -19 | -25 | -22 | 2 | 2.1 |
| 104 | 5 | -10 | -2.5 | 3 | 3.4 |
| 105 | -11 | -12 | -11.5 | 4 | 3.8 |
| 106 | -10 | -15 | -12.5 | 5 | 4.0 |
| 107 | -16 | -25 | -20.5 | 5 | 5.3 |
| 108 | -3 | -5 | -4 | 2 | 1.9 |
| 109 | -10 | -12 | -11 | 2 | 2.8 |
| 110 | -1 | -4 | -2.5 | 2 | 1.9 |
| 111 | -8 | -9 | -8.5 | 0 | 0.0 |
| 112 | 12 | 6 | 9 | 0 | 0.0 |
| 113 | -4 | -7 | -5.5 | 1 | 1.3 |
| 114 | 3 | -2 | 0.5 | 20 | 19.8 |
| 115 | -3 | -4 | -3.5 | 2 | 2.8 |
| 116 | 1 | -1 | 0 | 0 | 0.0 |
| 117 | 0 | 0 | 0 | 0 | 0.0 |
| 118 | -17 | -22 | -19.5 | 4 | 4.7 |
| 119 | -15 | -15 | -15 | 4 | 4.3 |
| 120 | -5 | -6 | -5.5 | 0 | 0.0 |
| 121 | -19 | -28 | -23.5 | 5 | 6.2 |
| 122 | 9 | -4 | 2.5 | 1 | 1.1 |
| 123 | -11 | -19 | -15 | 24 | 24.5 |
| 124 | 100 | 100 | 100 | 40 | 45.5 |
| 125 | -15 | -15 | -15 | 4 | 3.1 |
| 126 | -40 | -42 | -41 | 15 | 17.9 |
| 127 | 13 | 1 | 7 | 1 | 1.2 |
| 128 | -38 | -42 | -40 | 6 | 6.3 |
| 129 | -4 | -22 | -13 | 11 | 12.0 |
| 130 | 2 | -4 | -1 | 1 | 1.3 |
| 131 | -6 | -10 | -8 | 2 | 2.2 |
| 132 | -19 | -41 | -30 | 12 | 10.3 |
|  |  |  |  |  |  |
| Mean | -1.62 | -7.43 | -4.52 | 5.50 | 5.87 |
| Median | -4.00 | -8.00 | -5.50 | 2.50 | 2.80 |

**Tables 4a-4d: Thresholds (Based on Logit Regressions)**

Table 4a: Thresholds - Treatment 1 FixedNotKnown

| subject | threshold |
|---------|-----------|
| 1 | -8 |
| 2 | -4 |
| 3 | -100 |
| 4 | -10 |
| 5 | -8 |
| 7 | -33 |
| 8 | -7 |
| 9 | 18 |
| 10 | -49 |
| 11 | -6 |
| 12 | -38 |
| 13 | -8 |
| 14 | -12 |
| 15 | 15 |
| 16 | 3 |
| 17 | -49 |
| 18 | -11 |
| 19 | -16 |
| 20 | -41 |
| 21 | -1 |
| 22 | -100 |
| 23 | 1 |
| 24 | -12 |
| 25 | -21 |
| 26 | 5 |
| 27 | -2 |
| 28 | -14 |
| 29 | -1 |
| 30 | -24 |
|  |  |
| mean | -18.35 |
| median | -9.80 |

Table 4b: Thresholds - Treatment 2 RandomKnown

| subject | threshold |
|---------|-----------|
| 31 | -21 |
| 32 | -32 |
| 33 | -30 |
| 34 | -1 |
| 35 | -17 |
| 36 | -17 |
| 37 | -2 |
| 38 | -1 |
| 39 | -42 |
| 40 | -1 |
| 41 | -19 |
| 42 | 0 |
| 43 | -20 |
| 44 | -1 |
| 45 | -11 |
| 46 | 1 |
| 47 | -40 |
| 48 | 41 |
| 49 | -2 |
| 50 | -11 |
| 51 | -44 |
| 52 | -15 |
| 53 | -30 |
| 54 | -40 |
| 55 | -8 |
| 56 | -20 |
| 57 | 2 |
| 58 | -11 |
| | |
| mean | -13.94 |
| median | -13.19 |

Table 4c: Thresholds - Treatment 3 FixedKnown

| subject | threshold |
|---------|-----------|
| 59 | 1 |
| 60 | -100 |
| 61 | -64 |
| 62 | -3 |
| 63 | -23 |
| 64 | -6 |
| 65 | -19 |
| 66 | -85 |
| 67 | -6 |
| 68 | -15 |
| 69 | 16 |
| 70 | -2 |
| 71 | -2 |
| 72 | 1 |
| 73 | -3 |
| 74 | -21 |
| 75 | -100 |
| 76 | -9 |
| 77 | -42 |
| 78 | -14 |
| 79 | -8 |
| 80 | -7 |
| 81 | -19 |
| 82 | -37 |
| 83 | 2 |
| 84 | -75 |
| 85 | -6 |
| 86 | 2 |
| 87 | -17 |
| 88 | 0 |
| 89 | -1 |
| 90 | -6 |
|  |  |
| mean | -20.89 |
| median | -7.94 |

Table 4d: Thresholds - Treatment 4 RandomNotKnown

| subject | threshold |
|---|---|
| 91 | -27 |
| 92 | -4 |
| 93 | -18 |
| 94 | -26 |
| 95 | -2 |
| 96 | -1 |
| 97 | -1 |
| 98 | 7 |
| 99 | 100 |
| 100 | -25 |
| 101 | 1 |
| 102 | -3 |
| 103 | -16 |
| 104 | 0 |
| 105 | -11 |
| 106 | -13 |
| 107 | -23 |
| 108 | -5 |
| 109 | -11 |
| 110 | -7 |
| 111 | -10 |
| 112 | 5 |
| 113 | -7 |
| 114 | -10 |
| 115 | -10 |
| 116 | -2 |
| 117 | -1 |
| 118 | -24 |
| 119 | -2 |
| 120 | -7 |
| 121 | -26 |
| 122 | -2 |
| 123 | 27 |
| 124 | 100 |
| 125 | -19 |
| 126 | -23 |
| 127 | 1 |
| 128 | -33 |
| 129 | -5 |
| 130 | 0 |
| 131 | -9 |
| 132 | -32 |
| | |
| mean | -4.17 |
| median | -7.00 |

**Table 5a**: For each treatment, the number of subject with significant coefficients (at 5% level) for the logit regression where $\beta_1, \beta_2$, and $\beta_3$ are the coefficients of own payoff, opponent's payoff and kindness index, respectively.

| | FixedNotKnown | FixedKnown | RandomNotKnown | RandomKnown |
|---|---|---|---|---|
| only $\beta_1 > 0$ | 15 | 14 | 19 | 12 |
| $\beta_1 > 0$ and $\beta_2 > 0$ | 5 | 5 | 3 | 4 |
| $\beta_1 > 0$ and $\beta_2 < 0$ | 1 | 0 | 0 | 0 |
| $\beta_1 > 0$ and $\beta_3 > 0$ | 0 | 1 | 0 | 5 |
| $\beta_2 < 0$ and $\beta_3 > 0$ | 0 | 1 | 0 | 0 |
| $\beta_1 > 0, \beta_2 > 0$ and $\beta_3 > 0$ | 0 | 0 | 2 | 0 |
| $\beta_1 > 0, \beta_2 < 0$ and $\beta_3 > 0$ | 0 | 0 | 0 | 1 |

**Table 5b**: Logit regressions of accepting a negative offer for each treatment

| | FixedNotKnown | FixedKnown | RandomNotKnown | RandomKnown |
|---|---|---|---|---|
| own offer | 0.126*** | 0.075*** | 0.085*** | 0.63*** |
| | (0.011) | (0.007) | (0.008) | (0.006) |
| partner's offer | −0.003 | 0.071 | 0.018 | 0.042 |
| | (0.030) | (0.037) | (0.025) | (0.024) |
| constant | 0.850 | −5.618 | 1.218 | −2.331 |
| | (1.829) | (2.991) | (1.449) | (1.732) |
| # of observations | 1139 | 1024 | 1290 | 973 |
| Standard Errors are in parantheses. ***: significant at 1%, **:significant at 5% | | | | |

**Table 6a: Difference in Pair Thresholds:**

**FixedNotKnown Treatment**

Format: Difference   (Threshold1, Threshold 2)

| Pair | All rounds | Rounds 1-5 | Rounds 6-10 |
|------|------------|------------|-------------|
| **1-4** | 5    (-10,-5) | 3    (-8,-5) | 1    (-10,-9) |
| **2-18** | 13 (-9, -22) | 15 (-10, -25) | 3    (-7, -10) |
| **5-13** | 1    (-6,-7) | 10  (-15,-5) | 3    (-6,-9) |
| **8-19** | 6    (-8, -14) | 9    (-5,-14) | 3    (-10,-13) |
| **11-16** | 1    (0,-1) | 11 (-12, -1) | 0    (+2,+2) |
| **21-27** | 1    (-1, 0) | 9    (-1, +8) | 0    (-1, -1) |
| **23-26** | 4    (-4, 0) | 13 (+9, -4) | 2    (+2, +4) |
| **24-28** | 4    (-8, -12) | 10 (-7, -17) | 2    (-11, -9) |
| **25-29** | 22  (-19, +3) | 23 (-19, +4) | 1    (-2, -3) |
| **Mean** | 6.3 | 11.4 | 1.7 |

**Table 6b: Difference in Pair Thresholds**

**FixedKnown**

| Pair | All rounds | Rounds 1-5 | Rounds 6-10 |
|------|------------|------------|-------------|
| **59-80** | 1    (-1 0) | 9    (+8, -1) | 4    (-4, 0) |
| **65-71** | 22  (-23, -1) | 23 (-24, -1) | 12 (-14, -2) |
| **67-73** | 4    (-3,+1) | 6    (-3, +3) | 1    (-3, -2) |
| **72-79** | 7    (0, -7) | 1    (0, -1) | 9    (+2, -7) |
| **83-88** | 5    (1, -4) | 10 (-1, -11) | 2    (+1, +3) |
| **85-87** | 1    (-4,-3) | 8    (-4, -12) | 3    (+7, +4) |
| **86-90** | 1    (+5, +4) | 1    (+4, +3) | 5    (-8, -2) |
| **Mean** | 5.9 | 8.6 | 5.1 |

# 8   Appendix B

**Proof of Proposition (Equilibrium (Instrumental) Reciprocity):** We first prove that any optimal equilibrium must have the property that $x_i^C(w_i, w_j) = 1$ if and only if $w_i \geq -\ell_i$. Next we show that there is a unique optimal equilibrium with this property. Finally, in the section below, which demonstrates how to derive the optimal $\ell$, we demonstrate how $\ell$ varies with the Nash threat assumed in the punishment phase.

Consider two points in the second quadrant (that is, where $w_1 > 0$ and $w_2 < 0$): $A = (w_1^A, w_2^A)$ and $B = (w_1^B, w_2^B)$. Suppose that $w_2^A > w_2^B$, $x_2(w_1^A, w_2^A) = 0$ and $x_2(w_1^B, w_2^B) = 1$. In other words, player 2 approves proposal $B$ but vetoes proposal $A$, even though proposal $A$ gives player 2 a higher payoff. If this were an equilibrium, then player 2's no-deviation constraint must be met at point $B$. But then it must also be met at point $A$. It follows that, by choosing $x_2(w_1^A, w_2^A) = 1$ instead, we get an alternative equilibrium with a higher sum of joint payoffs—a contradiction.

The above argument implies that players' strategies must take the form $x_i^C(w_i, w_j) = 1$ if and only if

$w_i \geq -\ell_i$. It also implies that the no-deviation constraint, $w_i + \frac{\delta}{1-\delta} E_i \geq \frac{\delta}{1-\delta} N_i$, is exactly binding when $w_i = -\ell_i$:

$$-\ell_i + \frac{\delta}{1-\delta} E_i = \frac{\delta}{1-\delta} N_i, \tag{1}$$

Finally, it also implies that equilibrium payoff for player $i$ is given by

$$E_i \equiv \int_{\substack{w_i \geq -\ell_i \\ w_j \geq -\ell_j}} w_i\, f(w)\, dw - \int_{\substack{w_i \geq 0 \\ w_j \geq 0}} w_i\, f(w)\, dw.$$

Notice that $E_i$ is increasing in $\ell_j$ and decreasing in $\ell_i$.

We now show that there exists a unique efficient equilibrium, that is, one that maximizes joint payoffs. Suppose there were two such equilibria, corresponding to threshold levels $(\ell_i', \ell_j')$ and $(\ell_i'', \ell_j'')$ and leading to equilibrium payoffs $(E_i', E_j')$ and $(E_i'', E_j'')$, respectively. Without loss of generality, assume $E_i'' \geq E_i'$ and $E_j' \geq E_j''$.

Equation (1) and $E_i'' \geq E_i'$ imply $\ell_i' \leq \ell_i''$. By a similar argument, $\ell_j' \geq \ell_j''$. Since $E_i$ is increasing in $\ell_j$ and decreasing in $\ell_i$, this implies that $E_i'' \leq E_i'$. Given our starting assumption that $E_i'' \geq E_i'$, we conclude that $E_i'' = E_i'$, and so $\ell_i' = \ell_i''$. By a similar argument, we also conclude that $E_l'' = E_l'$ and $\ell_j' = \ell_j''$. ∎

$\diamond$ Derivation of equilibrium $\ell$: First we compute the value of $\pi^N$, equilibrium payoff in the static Nash game. Recall that there are two types of Nash equilibrium. In the weakly dominant strategy one (a player accepts a proposal if and only if his payoff is positive):

The area of the region where $w_i \geq 0$, for both $i$, is given by

$$\int_0^{100} x\,(100 - x)\, dx = \frac{500,000}{3}. \tag{2}$$

Straightforward calculations show that the total area of the set of proposals is given by 15,000. Since the distribution of $w$ is uniform over this set, it follows that $\pi^N$ is given by (2) divided by 15,000, or simply

$$\pi^N = \frac{100}{9}.$$

In the class of Nash equilibria in which any offer is rejected, the payoff of each player is equal to 0.

The next step is to compute the value of $\pi^E$, payoff along the repeated game efficient equilibrium path. The area of the shaded region in Figure 4 is given by

$$\int_{-\ell}^0 x\,(100 - (-x))\, dx + \int_0^\ell x\,(100 - x - (-x))\, dx + \int_\ell^{100} x\,(100 - x - (-\ell))\, dx,$$

48

or simply

$$\frac{2}{3}\,\ell^3 - \frac{1,000,000}{3} + \frac{1}{2}\,(100+\ell)(10,000-\ell^2). \tag{3}$$

It follows that $\pi^E$ is given by (3) divided by 15,000, or simply

$$\pi^E = \frac{\ell^3}{22500} - \frac{200}{9} + \frac{(100+\ell)(10000-\ell^2)}{30000}.$$

Given the values of $\pi^E$ and $\pi^N$, we can now derive the equilibrium value of $\ell$ by making the no-deviation inequality binding. We thus have

$$\ell + \delta\,\pi^E/(1-\delta) = 0 + \delta\,\pi^N/(1-\delta).$$

If the players use their weakly dominant strategy as a threat, first note that zero is a root. In fact, if $\ell = 0$, then $\pi^E = \pi^N$ and the no-deviation constraint holds trivially. Hence, we are left with a quadratic equation with the roots: $150 \pm 50\sqrt{36/\delta - 39}$, and it can easily be shown that only one of the roots (potentially) lies in the relevant interval, $[-100, 0]$. We thus have

$$\ell = 150 - 50\sqrt{36/\delta - 39}.$$

Solving for $\ell < 0$, we get $\delta > \frac{3}{4}$. Solving for $\ell > -100$, we get $\delta < .9$. So finally we have

$$\hat{\ell} = \begin{cases} 0 & \text{if} \quad \delta < .75 \\ -150 + 50\sqrt{36/\delta - 39} & \text{if} \quad .75 \le \delta \le .9 \\ -100 & \text{if} \quad \delta > .9 \end{cases}$$

In particular, $\delta = .8$ (parameter in the experiment) implies $\hat{\ell} = -27.53$.

If any offer is rejected as a punishment strategy, then $\ell + \delta\,\pi^E/(1-\delta) = 0$ and again only one root lies in $[-100, 0]$. When $\delta = .8$ implies $\hat{\ell} = 88.83$. Hence, any $\hat{\ell}$ between 27.53 and 88.83 can be sustained in the equilibrium. Since the sum of the offers are always positive, the efficient equilibrium is achieved when $\hat{\ell} = 88.83$.

**Proof of Proposition (Equilibrium with Altruistic Preferences):** Let $\Phi_i(w_{it}, w_{jt})$ be the altruistic utility function of the player $i$ such that $\frac{\partial\Phi}{\partial w_{it}} > 0$ and $\frac{\partial\Phi}{\partial w_{jt}} > 0$ for all $i, j = 1, 2$ Suppose on the contrary that the optimal trigger-strategy equilibrium exists in threshold strategies, in other words there

49

exists $\ell_i$ such that $x_i(w_i, w_j) = 1$ if and only if $w_i \geq -\ell_i$ for some $\ell_i < 100$. . In particular, the proposal $(-l_i, l_i)$ is accepted but player $i$ vetoes the proposal $(-l_i - \varepsilon, 100)$ for any $\varepsilon > 0$. Since $\frac{\partial \Phi}{\partial w_{jt}} > 0$, $\Phi_i(w_{it}, w_{jt})$ is continuous and $\ell_i < 100$, there exists an $\varepsilon$ as small as possible such that $\Phi_i(-l_i - \varepsilon, 100) > \Phi_i(-l_i, l_i)$. Hence, we get an alternative equilibrium in which both players would approve $(-l_i - \varepsilon, 100)$ and in this an alternative equilibrium with a higher sum of joint payoffs - a contradiction. ∎

This is an experiment in decision making. Money has been provided for this experiment by various research foundations. You will be paid for your participation and if you make good decisions you may be able to earn a substantial amount of money that will be paid to you when the experiment is over.

**The Experiment.**

You have been recruited to participate in this experiment along with a number of other people who are in the room with you. When the experiment starts you will be paired with one person in the room at random. This person will be your pair member for the rest of the experiment. The experiment will consist of 10 rounds with each round consisting of a random number of periods. While the number of periods will be random, there will always be at least 6 in any round. After period 6 is over in any given round, whether you proceed to the next period will be determined randomly as will be described below. So in any round you will play 6 periods for sure and maybe more.

When Round 1 starts you and your pair member will be shown a computer screen upon which two numbers will be shown, one indicating a potential payment to you and one a potential payment to your pair member. These payments are denominated in a fictitious experimental currency called francs which will be converted at the end of the experiment at a rate of 1 franc = 0.6 cents. The numbers will be drawn at random but all the pairs of numbers you will see will have the same two properties:

**1)Each payment to each pair member will be independently chosen from the interval [-100, +100].**

In other words, you will never see a number outside this range, and within this range each number will have an equally likely chance of being chosen, so there will be an equal chance that your number is –20 as it will be +85 as it will be –99, or +42 etc.

**2)The sum of the numbers must be positive and less than 100.**

This means that when the computer draws a number for you and your pair member and adds them up, if the sum is negative or greater than 100 the computer will throw that pair of numbers away and pick another. You will only see pairs whose sum is positive and less than 100.

When the payment pairs are drawn you will see a screen that says:

**Your franc payment _ _ _ _ _ _ _ _**

**Your pair member's franc payment _ _ _ _ _ _**

And you will be asked to approve or refuse the payments. If you approve you must click the approve button at the bottom of the screen; if you refuse, click the refuse button.

If both subjects approve the payment pair, you will both receive that amount indicated as a payment for the period. If either of you refuse, you both will receive nothing during that period.

When period 1 is over, you will be shown the results of that period by being informed of what your pair member chose (accept or refuse) and your payoff . We then proceed directly to period 2, which will be identical to period 1, that is, you will be shown a new randomly drawn payment pair and asked to accept or refuse. When period 2 is over you will be shown the results of that period and also your cumulative payoff up until that period and then proceed to period 3. This will happen for 6 periods. After the 6th period, the computer will randomly determine whether you move to period 7. It does this by flipping a coin that has a .80 chance of landing heads and a .20 chance of landing tails. If the coin lands heads, we proceed to period 7 and repeat the procedures above; and at the end the computer randomly determines whether we move to period 8. This will continue until the computer determines that you will not proceed to any more periods. When that is determined the computer will notify you that you will now be playing the last period in this round by announcing that **"This is the Last Period in this Round".** After you finish that period the current round of the experiment will end. In other words, you will be told when you are playing the last period in any round. When the round is over, you will be shown your payoffs for that round, proceed to the Round 2 and repeat the experiment identically. There will be 10 rounds in the experiment, but as you have seen, each round may have a different number of periods, depending on chance.

As stated above, all of the above payments will be denominated in a fictitious experimental currency called francs. At the end of the experiment, your payment will be converted into US dollars at the rate of 1 franc = 0.6 cents.

There is only one detail left to be explained: Within the first 6 periods of any round of the experiment, the periods we know we will play for sure, the number of francs you will receive when you accept a proposed payment pair will vary. More precisely we will multiply your displayed franc payoff by a "multiplier" depending upon the period the payoff is accepted (your pair member's franc payoff will also be multiplied by this number as well). The set of multipliers used is shown in Table 1 below:

**Period Multipliers**

| Period | Multiplier |
|--------|-----------|
| 1 | 3.05 |
| 2 | 2.44 |
| 3 | 1.95 |
| 4 | 1.56 |
| 5 | 1.25 |
| 6 | 1.00 |
| 7+ | 1.00 |

To illustrate what this table says, say that in period 1 of any round you and your pair member agree to a payoff pair that gives you a franc payoff of 20. In such a case instead of you being credited with 20 francs as your payment you would be credited with 20 x 3.05 = 61, where 3.05 is the multiplier associated with period 1 in the table above. If you agreed to the same payoff in period 5 you would be credited with 20 x 1.25 = 25, where 1.25 is the period-5 multiplier. Note that the multipliers decrease as we approach period 6, the last period you will engage in for sure, where the multiplier is equal to 1. It will remain equal to 1 for all succeeding periods; but, as we have explained above, in all succeeding periods the probability of continuing is equal to .8.

**Final Payoffs:**

Your final payoff in the experiment will be the sum of your earning over the 10 rounds of the experiments. That means that we will sum your franc payoffs earned in each round of the experiment and then convert them into U.S. dollars at the rate of 1 franc = 0.6 cents.