



UNIVERSITA' DEGLI STUDI DI TRENTO - DIPARTIMENTO DI ECONOMIA

CONFORMITY AND RECIPROCITY IN THE “EXCLUSION GAME”: AN EXPERIMENTAL INVESTIGATION

Lorenzo Sacconi
and
Marco Faillo

The Discussion Paper series provides a means for circulating preliminary research results by staff of or visitors to the Department. Its purpose is to stimulate discussion prior to the publication of papers.

Requests for copies of Discussion Papers and address changes should be sent to:

Dott. Stefano Comino
Dipartimento di Economia
Università degli Studi
Via Inama 5
38100 TRENTO ITALY

CONFORMITY AND RECIPROCITY IN THE “EXCLUSION GAME”: AN EXPERIMENTAL INVESTIGATION

Lorenzo Sacconi* and Marco Faillo**

(20/06/ 2005)

Abstract

Sacconi and Grimalda (2002, 2005a, 2005b) introduced a model in which two basic motives to action are understood as different type of preferences and represented by a *comprehensive* utility function: the first is *consequentialist* motivation, whereas the second is a *conditional willingness to conform with an ideal*, or a *moral principle*, which they call a *conformist*, or *ideal*, motive to action. A moral ideal is meant as a normative principle of evaluation for collective modes of behaviours which provides agents with a ranking of states of affairs resulting from strategic interaction expressing a greater or lesser consistency with the ideal. The principle moreover is seen as resulting from a (possibly hypothetical) contract between the agents involved in the interaction in an *ex-ante* phase. Thus, the normative principle boils down to a social welfare function that measures the consistency of outcomes with the normative prescriptions provided by the ideal. Hence, agents understand their own and any other agent's degree of conformity in terms of their contribution to carrying out the ideal given the others' expected action, and a person's own motivation to act in conformity with the principle increases with others' (expected) conformity. In other words, individual conformity with the principle is *conditional* on others' conformity with it, as perceived by the agent. This peculiar feature of *reciprocity* over others' behaviour calls for an extension of the usual equipment of decision theory, which is provided by the theory of Psychological Games (Geanakoplos *et al.*, 1989).

In this paper we design an experiment for preliminary exploration of the empirical validity of the conformist preferences model, applying it to a simple non cooperative game (the Exclusion Game) meant as the problem of dividing a sum between two active players and a third, dummy player (passive beneficiary). Results are encouraging. Behaviours dramatically change passing from the simple exclusion game to a three steps game, in which once the players have first played the typical non cooperative exclusion game, and before playing it again, they participate in a middle phase, where they anonymously agree on a principle of division. Having agreed on a principle, even though this agreement does not implies reputation effects nor is externally enforceable, induces a substantial part of players - who acted selfishly in the first step - to conform to the principle in the third phase. The additional condition being that they believe the other players will also conform to the agreed principle (what here does happen, as a matter of fact). These results strictly accord with the prediction of the conformist preferences model, but cannot be accounted for by alternative theories of reciprocity.

Keywords: conformist preferences, reciprocity, psychological games , fairness, experiments.

JEL classification: C7, C9,

* Department of economics, University of Trento and EconomEtica. lorenzo.sacconi@economia.unitn.it

** University of Trento and Sant'Anna School of Advanced Studies, Pisa. mfaillo@economia.unitn.it

The authors wish to thank Luigi Mittone, Gianluca Grimalda and Roberto Gabriele for useful suggestions and comments on the experimental design; Ivan Soraperra who produced the software; Marco Tecilla, Sara Martinatti, Manuela Seppi and the CEEL staff for their help in the conduct of the experiment.

1. Conformist preferences and the game of reference

1.1 Background theory.

Contrary to the economist's fixation with the model of rational economic man seen as uniquely *selfish* utility maximiser, a more reasonable and realistic understanding of rational action suggests that individuals weight different, possibly conflicting, motives to act. On the hypothesis that economic agents are motivated both by consequentialist (and mainly self-interested) and "conformist" preferences - that is the intrinsic motivation to act according to shared principles complied reciprocally by also the other interacting agents - Grimalda and Sacconi (2002, 2005a, 2005b) suggested that two classes of motives for choice can be grasped by two types of preferences of the Self and by their relative mathematical representation in the corresponding part of a *comprehensive* utility function. Thus, the utility maximisation model of rational economic man can be considerably revised, extending its explanatory and normative power at a substantial level. In this paper we describe the results of an experiment designed to test the explanatory capacity of Sacconi and Grimalda's model of choice. This first subsection summarises the theory because it furnishes also the hypotheses for the experimental study reported thereafter.

The model assumes what we call a *description relative* view point of motives to act. The same states of affairs generated by strategic interaction can be described in different ways according to their relevant characteristics. A first description of states views them as *consequences*. Consequences may be described as attributed only to the acting Self - what happens to the decision-maker in any state. This description is the basis of Self-interest: the Self defines preferences over consequences, which are Self-referred. By contrast, consequences may be attributed to any person (extended consequences) in so far as they can be understood as consisting in what happens to any whatever individual. To make sense of any impartial consequentialist ethics like utilitarianism, or altruism, consequences must be extended like that. In general, if a player defines his preferences only on states *described as consequences*, then he has *consequentialist personal preferences*.

The second type of preferences is what we call *conformist personal preferences*. Description of states is no less important here, but they are now described as sets of interdependent actions played by the participants in any generic game (call such an action profile σ), characterised in terms of whether or not they are consistent with a given abstract principle.

First, the principle can be captured by a function of individual first type utilities attached to states which measures the fairness of welfare distribution in each state of the world. A pattern of behaviours (a vector of strategies) is defined as perfectly *deontological* if it is fully consistent with the abstract principle of fairness - that is, if it maximises the function just defined. We call this state the *ideal*. Hence, the fairness function generates an ordering of states in terms of their consistency to the principle. Secondly, we then determine the individual degree of conformity with the ideal relative to each strategy choice of any player by seeing whether and how much the *ideal* comes about through that choice of the player, given what he believes about other parties' choices. Moreover the motivational weight of this conformity to the player himself depends on the expected reciprocal conformity of the other players given what they expect on the part of the first player. In fact, for each strategy combination, the intensity of each player's conformist preference will depend on a measure of expected reciprocal conformity, which is based on:

- (a) an index of the extent to which the player himself contributes to fulfilling the ideal by conforming with or deviating from it, given what he believes about the other player's choice;
- (b) an index of the extent to which the player believes that the other player contributes to fulfilling the ideal by conforming with or deviating from it, given what the second player believes (and the first player believes that the second player believes) that the first player will do.

It follows that conformist preferences rest upon a measure of reciprocal conformity with the ideal, which on the other hand depends on a principle-based ordering of states according to a deontological property they possess. There is an indubitable relation between the two kinds of preferences - consequentialist and conformist. In fact, in order to define fairness, we look at distributions of the payoffs deriving from the *first type of preference* - i.e. material utilities. But this does not reduce second-type preferences to first-type ones. First-type utilities are no more than rough materials for the criterion defining second-type preferences. What matters for the second description are not consequences or material payoffs as such, but a distributive property defined over payoffs, which is expressed by the function representing a principle of fairness.

Thus characterisation of second-type preferences accords more with deontology than consequentialism. The more an expected state of affairs is consistent with the ideal, the more an individual action belonging to it is preferred by a player (granted that the player

contributes to generating the state and he expects that the other player reciprocally contribute in the same way). Moreover, there is no reason to link deontology with a belief that there is some objective source of value that has ontological reality “out there” (independently of the decision maker’s affections). In fact, while conformist preferences depend on degrees of deontology, deontology itself may nonetheless be understood, as it is here, simply as individual compliance with a fair distribution principle *that players could have rationally agreed in an ex ante hypothetical bargaining situation*.

Once these concepts have been translated into a formal model, a player’s *comprehensive utility function* consisting of two parts (which we assume to be separable), i.e. the representations of consequentialist and conformist preferences, can be defined as

$$(1) \quad V_i(\sigma) = U_i(\sigma) + \lambda_i F[T(\sigma)]$$

where U_i is the traditional player’s i “consequentialist” utility for state σ (a given strategy combination of any generic game), the weight λ_i - which may be any positive real number - is an exogenous psychological parameter that expresses how important the conformist component is within the motivational system of player i (we could call it the player’s i ‘maximum disposition to act according to conformist reasons’, granted that certain conditions do apply), and F is a function representing reciprocal conformity with principle T which in turn is a function taking a value for each state σ .

First required is specification of a form of the fairness-function T which formally represents the ideal. This must be a mapping from the set of states (and first-order utilities attached to them) to a fairness ordering ranging over states. A characterisation in contractarian terms of the ideal principle T is given by the Nash bargaining solution, i.e. the *Nash social welfare function* N (or *Nash product*)

$$(2) \quad T(\sigma) = N(U_1, \dots, U_N) = \prod_{i=1}^N (U_i - c_i)$$

where c_i represents the reservation utility that agents can obtain when the process of bargaining breaks down. Second, as far as the function F is concerned, defined are two personal indexes of conformity, which will be compounded in a measure of mutual expected

conformity. It will and enter the utility functions of the players in so far as it will influence the weight of the conformist motivation within the individuals system of preferences.¹ In this construction we take the point of view of player i (any other player j 's perspective is symmetrical).

Player i 's personal index of conformity:

This is player i 's degree of deviation from the ideal principle T (which varies from 0 to -1), due to player i 's choice, given her expectation about player j 's behaviour. It is normalised by the magnitude of difference between player i 's full conformity and no conformity at all conditional on player j 's choice

$$(3) \quad f_i(\sigma_i, b_i^1) = \frac{T(\sigma_i, b_i^1) - T^{MAX}(b_i^1)}{T^{MAX}(b_i^1) - T^{MIN}(b_i^1)}$$

where b_i^1 is player i 's belief concerning player j 's action, $T^{MAX}(b_i^1)$ is the maximum attainable by the function T given i 's belief, $T^{MIN}(b_i^1)$ is the minimum attainable by the function T given i 's belief, $T(\sigma_i, b_i^1)$ is the effective level attained by T when the player adopts strategy σ_i , given his belief about the other player's behaviour.

Estimation function of the second player's index of conformity with the ideal:

This is player j 's degree of deviation from the ideal principle T (which also varies from 0 to -1), as seen through player i 's beliefs - also normalised by the magnitude of difference between player j 's full conformity and no conformity at all, given what he believes (and player i believes that he believes) about player i 's choice

$$(4) \quad \tilde{f}_j(b_i^1, b_i^2) = \frac{T(b_i^1, b_i^2) - T^{MAX}(b_i^2)}{T^{MAX}(b_i^2) - T^{MIN}(b_i^2)}$$

where b_i^1 is player i 's *first order* belief about player j 's action (i.e. formally identical to a strategy of player j), b_i^2 is player i 's *second order* belief about player j 's belief about the action adopted by player i (i.e. formally identical to a player i strategy predicted by player j).

¹ Grimalda and Sacconi (2002) elaborate on Rabin (1993) in order to define the model of reciprocity.

These indexes are compounded to construct the following conformist component of the utility function

$$(5) \quad \lambda_i \left[1 + \tilde{f}_j (b_i^2, b_i^1) \right] \left[1 + f_i (\sigma_i, b_i^1) \right]$$

From this formula we may state the following: if player i perfectly conforms with the ideal, given her expectation, while the player j is also expected to perfectly conform, then the two individual indexes take zero values, so that the resulting utility value due to conformism is λ_i . By contrast, if a player does not entirely conform, while not expecting the other player entirely to conform either, then the two indexes take negative values (possibly -1). Thus the utility calculation for conformist reasons reduces to $(1-x)(1-y)$ (possibly both equal to zero) times the weight λ_i and yields less than λ_i (possibly zero) as the conformist utility value.

The comprehensive utility function V_i , hence, takes the form of the linear combination of the two components, with reference to each state described in terms of both the individual strategy choice and the individual beliefs system over the other player's strategy choice (note that in equilibrium beliefs meet choices, so that a belief accurately describes what the other player does)

$$(6) \quad V_i (\sigma_i, b_i^1, b_i^2) = U_i (\sigma_i, b_i^1) + \lambda_i \left[1 + \tilde{f}_j (b_i^2, b_i^1) \right] \left[1 + f_i (\sigma_i, b_i^1) \right]$$

This suggests that if a player predicts reciprocal conformism (so that conformist motivation effectively enters the utility function with weight λ_i), as long as the weight λ_i is high enough, it is possible that the overall utility function reverses the preference for a strategy choice σ_i with respect to the same player i 's simple consequentialist preferences represented by $U_i(\sigma_i, b_i)$. For example, it may induce the players to select strategies that they would never have chosen if they had relied on their material utility only.

1.2. A reference situation: the Exclusion Game

Conformist preferences account for ideal-driven choices based on the motivational force of reciprocal conformity with impartial and impersonal abstract principles of fairness. Hence it seems natural to use them in an attempt to explain behaviours such that a decision of fairly including a weak party in the sharing of benefits is taken by a group of strong players. The stylized situations of reference is one where "inclusion" is due not to the attitudes (or

intentions) expressed through effective interaction by the weak toward the strong parties (kindness, friendship, spontaneous cooperation and the like, which may affect their welfare position), but simply is a decision taken out of a sense of commitment. In other words strong players recognise an intrinsic value in reciprocating (amongst themselves) conformity with a principle of fairness. The typical situation is a social interaction between strong players and a weak player such that their mutual interaction makes a social surplus affordable, but only the strong players have decision influence over the allocation and distribution of the social surplus, whereas the weak party has neither a voice in this decision nor any retaliation threat at her disposal. Strong players can then decide to include the weak player in the fair sharing of the surplus, or alternatively to exclude the weak player and share out the pie among themselves. This is why we call this situation the “Exclusion Game”. Of course, one should not forget the role of reciprocity. Nevertheless reciprocity in this case cannot have anything to do with “intentions” showed by the weak party, which in fact is so weak that it can be represented as a dummy player, i.e. a player with no move in the game, and hence incapable to show intention through actions. On the contrary reciprocity in conforming with a fairness principle works through the preferences of the group of players with decision influence and concerns their reciprocal conformity with the fairness ideal. Explanation of their inclusive behaviour does not involve any kindness, reciprocity, direct interaction or attitude whatsoever shown by the weak player toward the strong ones. In short, strong players’ inclusive behaviour is born out of a ‘sense of commitment’, it amounts to ‘noblesse oblige’ .

In order to give a formal description of the exclusion game, we consider a situation (a non-cooperative game) in which two individuals (player 1 and player 2) must decide how to allocate a sum of money R among themselves and a third individual (player 3), who does not have an active role in the allocation decision but whose payoff is determined by the choices made by the two other players (active players). In particular, active players can choose between two alternative strategies: first, asking for a large share of R , i.e. high demand, $d_i^h = \frac{R}{2}$, which jointly amount to the whole pie; or second, asking for a small share of R , compatible with a fair distribution of the surplus, i.e. the low demand $d_i^l = \frac{R}{3} = s$, with $i = \{1,2\}$. The third player’s payoff is the remaining share of R after the two demands of the

active players have been met, i.e. $R - (d_1 + d_2)$. For example, in the case of low demand strategies by both the active players, the third player's payoff is $s = R - (d_1^l + d_2^l) = R/3$

Figure 1. The Exclusion Game. Payoff matrix

	d_2^l	d_2^h
d_1^l	d_1^l, d_2^l, s	$d_1^l, d_2^h, s/2$
d_1^h	$d_1^h, d_2^l, s/2$	$d_1^h, d_2^h, 0$

As shown in figure 1, if both the active players decide to ask for half of the total sum R (high demand strategies), the third player's payoff is zero; if one of the two active players decides to ask for only one third of R, while the other one chooses to ask for half of R, third player's payoff is R/6. An equal division of R among all the players - including the dummy third player - (R/3 each), results when players 1 and 2 ask for only R/3.

Given the assumption that player 3 is *dummy*, as far as we maintain that the players are motivated only by the intent to maximise their material payoffs, the only equilibrium in dominant strategies of this game is the one in which both active players choose to ask for $d^h = \frac{R}{2}$. Thus, the Exclusion Game played by self-interested players will induce effective exclusion of the player with no influence on the allocation decision. This is only a preparatory move, however. The next section illustrates how matters changes when conformist preferences and reciprocity among the active players enter the picture..

1.3 The exclusion game under conformist preferences and reciprocity.

In order to apply the theory of conformist preferences to the Exclusion Game, we must first re-describe states of affairs resulting from the game in terms of their consistency to the ideal of fairness formalised as the Nash bargaining solution (or Nash product) (where it is assumed that the status quo is zero for each player). This we assume to be the principle of fairness agreed upon by all the players (active or otherwise) in an hypothetical (i.e. *not externally enforceable*) agreement.

On applying the Nash product to the states resulting from the possible plays of our game we obtain the following fairness ordering of the four strategy combinations

$$(7) \quad d_1^l d_2^l > d_1^h d_2^l = d_1^l d_2^h > d_1^h d_2^h$$

based on

$$(8) \quad \begin{aligned} T^{MAX}(d_1^l, d_2^l) &= N(d_1^l, d_2^l) = d_1^l d_2^l s \\ T(d^h, d^l) &= N(d_1^h, d_2^l) = d_1^h d_2^l \frac{S}{2} = N(d_1^l, d_2^h) = d_1^l d_2^h \frac{S}{2} R \\ T^{MIN}(d^h, d^h) &= N(d_1^h, d_2^h) = d_1^h d_2^h 0 = 0 \end{aligned}$$

We can use these values to compute the overall utility values on the basis of the conformity indexes for each pair of actions and for the relative beliefs. In this new context the appropriate notion of equilibrium is that of *Psychological Nash Equilibrium* (Geanakoplos et al. (1989)²), which is an extension of the Nash equilibrium for situations in which expectations enter the player's utility function.

Accordingly, given the players utilities defined as function of their beliefs, we can easily compute the psychological equilibria of the game played by agents with conformist preferences. Strategy combinations that were not Nash equilibria in the basic game can now be defined as psychological equilibria. In particular, (d_1^l, d_2^l) is a psychological equilibrium *once it is granted that the weights λ_i are sufficiently high*:

$$(9) \quad V_1(d_1^l, b_1^1 = d_2^l, b_1^2 = d_1^l) > V_1(\sigma_1, b_1^1 = d_2^l, b_1^2 = d_1^l) \Leftrightarrow \lambda_1 > d_1^h - d_1^l$$

$$(10) \quad V_2(d_2^l, b_2^1 = d_1^l, b_2^2 = d_2^l) > V_2(\sigma_2, b_2^1 = d_1^l, b_2^2 = d_2^l) \Leftrightarrow \lambda_2 > d_2^h - d_2^l$$

Put otherwise: given the vector of strategies (d_1^l, d_2^l) , every player i 's overall utility from strategy d_i^l - assuming a system of mutually consistent beliefs according to which each player predicts with probability 1 the symmetric strategy d^l by the opponent - is greater than the overall utility gained by deviating to the alternative strategy d_i^h , and this holds simultaneously true for both the active players. In our example this condition is satisfied when the weight of conformist preferences λ_i compensates for the loss of material utility deriving

² See Appendix 1 for details.

from the decision to comply with the ideal. Under these conditions there exists a psychological equilibrium of the game such that players 1 and 2 choose to ask for the lowest share of the total sum that guarantees an equal distribution of R among all three players. Thus, one of the equilibrium solutions of the psychological Exclusion Game is the effective inclusion of the third inactive party in the sharing of the surplus. In the event that players have strong preferences for reciprocal conformity with the hypothetical social contract ideal of fairness, and if they have consistent reciprocal beliefs in that regard, a solution may be inclusion, not exclusion.

Note, however, that this strategy combination is not the only equilibrium of the game: also (d_1^h, d_2^h) is a psychological equilibrium when a system of beliefs exists such that both player 1 and 2 predict with probability 1 that nobody will conform with the principle, in that they have higher-order beliefs coherent with this expectation. In particular, if player 1 believes that the opponent will choose the worst action with regard to the moral principle (first order belief), and if he also believes that player 2 believes that 1 will choose the same action (second order belief), neither the opponent nor player 1 have incentives to respect the moral principle by acting against their material self-interest.

In the following sections, after having introduced the experimental design and procedure (sec. 2), we will show (sec.3) how this model can be used to formulate predictions about the choices made by subjects involved in the experiment.

2. Experimental design and procedure

The experiment took place at the *Computable and Experimental Economics Laboratory (CEEL)* of the University of Trento and it consisted of six sessions with 15 subjects, for a total of 90 participants.³ Each subject received a show-up fee of €5 for participation.

Each session was divided into three phases and lasted one hour on average.

In phase one the subjects played a version of the Exclusion Game. They were assigned to groups composed of three members. Within each group, subjects were randomly attributed the roles of G1 and G2 and G3. G1 and G2 were invited to play a game in which they had to decide how to allocate a sum of money ($S = €12$) between themselves and the third player, who did not have any active role in the game. In particular, active players were able to decide

³ Participants were all students at the University of Trento (mainly from economics, law and sociology courses), recruited by responding to ads posted at the various departments.

how much of the sum to ask for themselves (d_1, d_2), selecting one of three possible strategies: 25%, 33% or 50% of S . Active players' payoffs corresponded to d_1 and d_2 , while the third player's payoff was $S-(d_1+d_2)$ (Figure 2).

The subjects played the game three times, in three different rounds. At the beginning of each round the three roles were randomly assigned to the members of the group. The selection mechanism was designed so that each player was able to take each of the three roles G1, G2 and G3 in turn. The subjects were told that at the end of the experiment the software would extract one of these three rounds at random, and the player's earning for phase 1 would be determined according to the outcome of that round.⁴ The game was played anonymously and subjects were not aware of the previous rounds' outcome. This procedure produced two observations for each player in this phase: his choice in the G1 role and his choice in the G2 role. That is for each player we have two rounds of observable choice at phase one.

Figure 2. The experimental Exclusion Game. Payoff matrix

		G2		
		3 (25%)	4 (33%)	6 (50%)
G1	3(25%)	3,3, (6)	3, 4, (5)	3, 6, (3)
	4(33%)	4, 3, (5)	4, 4, (4)	4, 6, (2)
	6(50%)	6, 3, (3)	6, 4, (2)	6, 6, (0)

In phase two the subjects were assigned to new groups consisting of three anonymous members. Without definition of roles, they were invited to agree upon an hypothetical rule for the allocation of a sum between two active players and one non-active player by means of a voting procedure. The agreement was to be reached by repeatedly playing the voting procedure until unanimity was reached, within a given limit of trials. No explicit communication was allowed among the players of any given group. In particular, after they were informed that in the following phase they would play a game like the one played in the

⁴ This is an application of the procedure known as *random lottery incentive system* (Starmer and Sugden (1991) and Cubitt, Starmer, and Sugden (1998)). Adopting this procedure, the round in which the subject has played in the G1 or G2 role is selected with a probability of 2/3, which is the same probability of being extracted as G1 or G2 in a one-shot version of the game. Note that, if we look at the third phase of the experiment, using this mechanism we can always compare the choice of each of the players that in that phase have an active role with his choice in the first phase.

first phase, they were requested to vote for one of two general principles and one among some more specific rules deduced from the selected general principle (Figure 3). Subjects were told that groups which reached unanimous agreement by voting for the same principle within five trials would pass to the voting on the specific allocation rule, upon which the groups had to agree within ten trials. A lack of unanimity after the last of the trials would prevent subjects from entering the third phase.

At the beginning of this phase the experimenters informed the subjects about the voting procedure, stressing the correspondence between the specific rules and the game strategies of phases one and three. Absolute anonymity was guaranteed and the subjects were not allowed to communicate throughout the procedure.⁵

Figure 3: Second phase. Principles and rules

	PRINCIPLE 1.				PRINCIPLE 2.			
Principles	Every player should share the benefits, in particular, who has not the possibility to choose should not receive less than				People who play under a decisional role could claim a higher share of benefits.			
		↓				↓		
	G1	G2	G3		G1	G2	G3	
Rules	1.1	33%	25%	42%	2.1	50%	25%	17%
	1.2	25%	33%	42%	2.2	33%	50%	17%
	1.3	33%	33%	33%	2.3	50%	50%	0

In phase three, with the composition of the group unchanged, G1, G2 and G3 roles were randomly assigned to the members of each group that agreed upon a given principle.

The subjects were involved in the same game as in the first phase, but now active players had the additional option of choosing between implementing the rule that they had agreed in the second phase or choosing one of the alternative strategies. If a player decided to implement the rule, then the corresponding strategy would be automatically selected, otherwise the strategy would be removed from his strategy set. Thus, for example, if player *i* was part of a group that in the second phase had reached agreement on rule 1.2 in figure 3 and if in phase three, playing the role of G1, he decided to implement that rule, then strategy ‘4’ was automatically selected.

⁵ See Appendix 2 for a detailed description of the voting procedure.

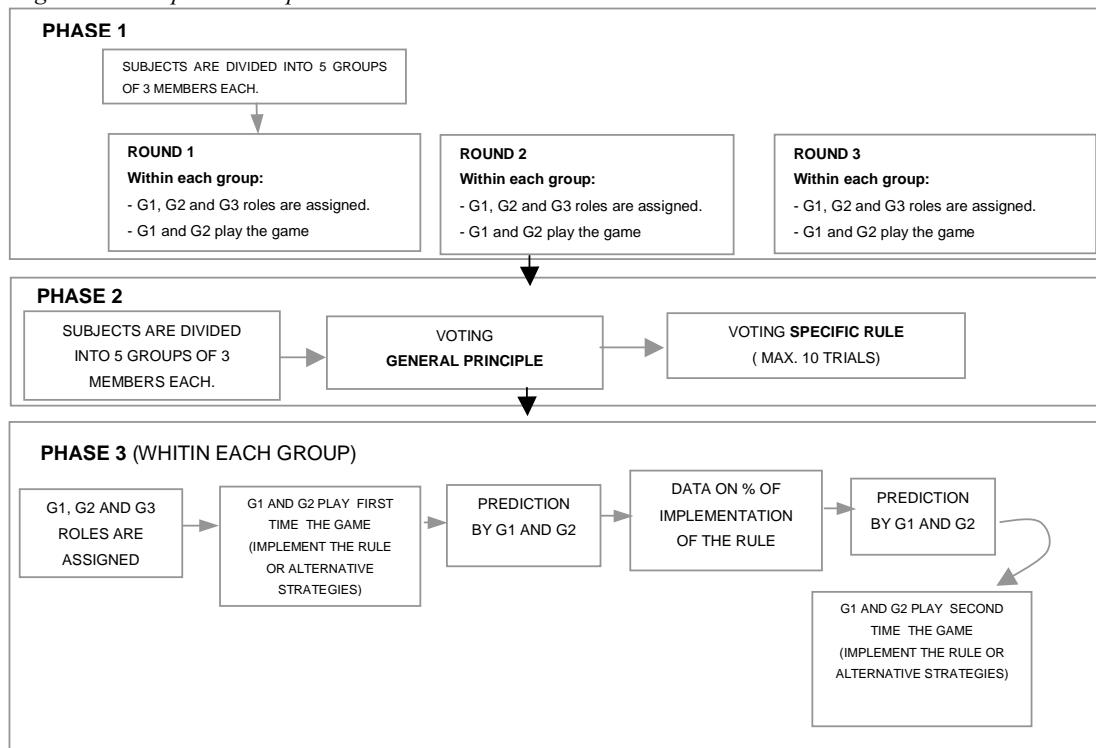
Just after their choice, active players were asked to express their expectation about the opponent's behaviour by guessing the outcome of the game.⁶ Data on the number of subjects that chose to implement the rule were collected and communicated to active players, who were asked to play again, in the same roles, using the same procedure as adopted at the beginning of the third phase, without knowing the outcome of the previous game.

At the end of the third phase, earnings were computed and players were paid. The total earning corresponded to the payoff of the extracted round of the first phase plus the payoff of the two games played in the third phase.

Then subjects were asked to fill out a short questionnaire on their explanation for doing certain decisions during the experiments. (see appendix 4).

A scheme of the experimental procedure is given in figure 4. The correspondence between the structure of the theory under control and the experimental design should be clear. In the first phase players enter the game without being able to agree upon a general principle of division.

Figure 4. The experimental procedure



⁶ We asked the player to indicate the cell of the payoff matrix in which he thought the game would end. In this way we avoided explicitly asking for his opinion about the opponent's willingness to conform with the rule.

The second phase corresponds to the “constitutional” step. Here the players can agree on both a general principle and a specific rule for the division of the sum. Note that one of the rules, the (33%, 33%, 33%) one, corresponds to the outcome of the game in which the Nash welfare function is maximised.

Playing this phase does not produce payoffs as such, because it only involves agreement on the rule that may be used to resolve a game later on, when the players will again have complete freedom of choice relative to the strategy to be implemented. However failure in the constitutional process is costly, because it prevents players from accessing the last phase, in which they can earn an additional sum of money.

The choice between the two principles should not be interpreted as simply a choice between moral conduct and pure amoral, and intrinsically censurable, self-interested behaviour. In fact, the inclusionist and exclusionist principles have been presented as if they express two alternative, but nevertheless admissible, positions. In the case of inclusion, the allocation of the benefits is assumed to be independent of the active role played by the subjects in the game (i.e. it seems as morally inappropriate as a basis for judgment whether the player holds the power to exclude or not); in the case of exclusion, the principle of equal distribution takes as relevant the active roles in the interaction and hence it is applied only within the active players set.⁷

In the third phase, the same players, after they have reached agreement on a division rule, are again involved in the Exclusion Game for two times. First they can choose whether or not to implement the rule they agreed. They are completely free to make a decision contrary to their previous agreement on a rule. In fact, even though agreement has been reached, at this stage its benefits are over, and the subjects are asked to look forward to the exclusion game as a new phase of the experiment. Hence, if the subjects have purely consequentialistic orientation they will forget everything about the agreed principles or rules, and will

concentrate on the possible outcomes of the exclusion game understood simply as consequences (material payoff). Otherwise considering the agreed principles and rules (and beliefs about compliance by other players) clearly denotes a deontological orientation. Asking the players about their expectations concerning the outcome of the game enabled us to deduce the players' beliefs about the opponent's willingness to conform with the rule. Finally, by asking the subjects to play the game again, after they had been told the percentage of people who had decided to implement the rule, we wanted to test the effect of such information on the players' beliefs and choices.

3. Empirical predictions and results.

3.1 Empirical predictions

What should we expect the results of the experiment to be if we assume that the players had conformist preferences?

The answer to this question is provided by direct application of the model presented in section 1.1. We begin by applying the Nash bargaining function to the outcomes of the game used in the experiment (2). Figure 5 reports the 'fairness values' given by the function in correspondence to the various states resulting from playing each strategy combination of the game .

Figure 5. Application of model to the Exclusion game. Nash welfare function values.

	3	4	6
3	54	60	54
4	60	64	48
6	54	48	0

⁷ One could criticise our decision not to sharpen the moral distinction between the self-interested rule of conduct and the other-regarding rule of conduct. On the other hand, were the choice framed as one between a moral principle and an amoral behaviour at all some reader would have been right in suggesting that we could have induced the subjects' choice by means of some interiorised cultural norm, not letting them to be free in agreeing on the constitutional principle they see more fit for this situation.

To sum up, those who chose the 'only-the-active-players-have-a-legitimate-claim-to-a-share-of-the-pie' principle showed to adhere to a sort of 'moral egoism' based on the intuition that only having effective influence on the result legitimise a claim (something like 'to whichever according to his power'). On the other hand those who chose the 'equal division' principle realised that being one of the participants in the choice over the constitution taken under a veil of ignorance concerning whether a player will be active or dummy, which will allow, once agreed, to proceed to play the division game, legitimise an equal claim on the final distribution no matter the random selection (by a natural lottery) of the role the players will take later on.

Hence, from the conformity indexes attached to each outcome of the game, we can compute the individual comprehensive utility values, assuming that in each state the players' beliefs reciprocally predict exactly the strategy chosen by the opponent. These values are reported in Figure 6.

Figure 6. Application of the conformist preferences model to the Exclusion game. Payoff matrix.

	3	4	6
3	3,3	$3 + \frac{3}{4}\lambda_1, 4 + \frac{3}{2}\lambda_2$	3,6
4	$4 + \frac{3}{2}\lambda_1, 3 + \frac{3}{4}\lambda_2$	$4 + \lambda_1, 4 + \lambda_2$	4,6
6	6,3	6,4	6,6

As one can see, if $\lambda_1 > 6 - 4 = 2$ then player 1 will prefer strategy '4' to strategy '6', and the same holds for player 2. Thus the strategy combination (4, 4) is a psychological equilibrium if $\lambda_i > 2$ and if the players' reciprocal beliefs are coherent with these strategies. But if player 1 believes that player 2 will not choose the strategy that produces the outcome closest to the ideal one, he will do the same, choosing strategy '6'. Because the same holds for player 2, the strategy combination (6, 6) is a psychological equilibrium as well. It follows that together the choice of the principle that may enter their indexes of conformity influencing the motivating force of conformism, empirical prediction about the solution of the game will depend upon what we can say about the combination of the players' reciprocal beliefs and the absolute weight of the conformist disposition.

In phase 1, players have no information about the type of their opponents, nor can they refer to any pre-existing agreement about the way in which the game should be played. Thus, there is no basis for conformist preferences (there is no agreed principle to comply with nor any reason to expect compliance by the counterpart). Even though the players could in principle have a high level of λ_i , this weight simply remains inactive. We should therefore expect players, even those with high conformist disposition, to believe that their opponents will

choose the strategy that maximises their own self-material interest, and consequently will ask for €6 (50% of S).

***Prediction 1.** In phase one, the choices of players motivated by conformist preferences will not be different from the choices of self-interested players. We will consequently find that they choose strategy '6'.*

In the second phase players must choose a rule on how to play an hypothetical Exclusion Game that may be played at a later moment. They know that if they are able to agree upon some principle of division, they will be able to play the Exclusion Game later, even though they do not yet know in what role they will play it again. This is typically a constitutional perspective. Such a perspective allows for the choice of general principles and rules of behaviour, incorporating a view of fairness. According to a contractarian approach to the constitutional choice of principles, players will take an impartial perspective: that is, they will judge the outcomes of the game from the point of view of each of the three roles in turn, and then choose a principle and a rule acceptable from whichever point of view. This implies a solution that must be invariant to the permutation of the individual points of view, that is, equal distribution of the surplus - if it is available within the payoff set - given what is claimed as baseline by every player in the constitutional choice. Note that within the “cooperative payoff space” defined by the Exclusion Game, rational bargaining according to the Nash bargaining solution would select the ‘equal division’ outcome (in coherence to the ‘invariance to symmetries’ and the Pareto postulates, granted the *status quo* is zero).

In this setting, ‘equal distribution’ is also an intuitively obvious choice, i.e. one with high ‘salience’. Given that agreement in this phase is a necessary condition for accessing the third phase, players may vote for the most salient rule to co-ordinate their choices in a limited number of trials. Salience, of course, may depend on the simplicity of the symmetric distribution; on the other hand, cognitive simplicity may also be connected to the fairness of equal division. Whether the cognitive simplicity or the intuitive fairness of a symmetric distribution comes first is difficult to say. We are here tempted to say that the cognitive and ethical features of symmetry are quite interlocked.

However, we must remind that available to conformist players at this step is also an agreement over the ‘the powerful players get all the pie’ principle, which notwithstanding its

crudeness is a possible second principle of division. Hence, we cannot uniquely predict that conformist players will choose the equal division principle. Conformity enters the picture only once a principle has been chosen in the constitutional phase, whereas the nature of the specific principle chosen depend on the proper understanding of the contractarian nature of the constitutional phase, which is independent of conformity *per se*. (We might say that also a constitutional choice of the “egoist” principle, even though it may reflect a misunderstanding of the symmetry of the contractarian choice, once it were made and conformed to, would be consistent with the model of conformist preferences). Thus we can only predict from our normative model of the constitutional nature of the decision phase over principles, that the equal division principle will have some intuitive force. Note that this is also a methodological necessity if the experiment must be able to test in effective way the hypothesis of conformist preferences. In fact, should conformist players - which at the first step decided to act out of their simple self-interest - decide at the second step uniquely to agree over “the powerful take all the pie” principle, and then going on by deciding to comply with this principle, then no evidence of change in the players’ behaviours could be observed through the experiment, since self-interest and conformity would dictate the same behaviour. No falsification could be provided this way against the conformist hypothesis. On the other hand, if a significant share of players, who made a selfish choice at the first phase, subscribe to an equal division principle at the second phase, then we have a clear empirical basis against which conformist theory can be tested. We have simply to see whether the mere fact of a ‘constitutional’ agreement over rules - which opens the door to the opportunity to play again a beneficial division game at phase two - is capable to activate motivational forces that will drive players to conform with the agreed principle (granted that players entertain the appropriate beliefs), changing their conduct with respect to how they behaved in the first phase.⁸ Thus, our second prediction is crucial to the falsification power of our experiment.

Prediction 2. *In phase two, a rule that assigns equal payoffs to all the players will be chosen by a significant part of the participants.*

⁸ Note that this implies that only a subset of the observations consistent with the theory may have crucial discriminating force amongst different theoretical hypotheses over rational action, and we mainly are interested to produce exactly this kind of evidence. This would have justified us also in stressing a bit more the ethical nature of the second phase decision in order to test the level of conformism in the third phase.

Assume that the players have now reached phase three. Hence they must have been able to agree on the same principle and rule. If this is enough for them to believe that a chosen principle and rule will also be played by the other players who have agreed on the same principle and rule, then their reciprocity-based conformist preferences will be activated (both deviations from conformity indexes are close to zero), granted that exogenous weights attached to non material motivations are significant. Hence, a conformist player will comply with the principle and the rule. If we hypothesize that the exogenous weight of conformist motivation is a psychological feature widespread in the population, and granted that we predict that a significant part of the players will have chosen an equal division rule, then we must expect a significant number of them to choose strategy '4'. What is most important, however, is that, if the players are conformist, we must expect the largest part of those who agreed on the equal division principle and rule to comply with the rule in the third phase, if we have some evidence that they believe that the rule will be followed by the others. Moreover, if these players are told that most of the players implement the principle and the rule (in a way that backs the relevant belief), then, if they are truly conformist, they will be more motivated to play in accordance with the principle and the rule themselves.

Prediction 3. *Players with conformist preferences who (having agreed on a rule) predict an outcome of the game compatible with a belief about the opponent's willingness to implement the agreed rule, will implement the rule as well.*

Prediction 4. *Given predictions 2 and 3, a significant part of players (endowed with conformist preferences) will request '4'. Moreover most of the player (endowed with conformist preferences) who behaved according to prediction 2 and satisfy prediction 3, will request '4' in the third phase.*

Prediction 5. *Once conformist players learn that a large part of the players that have chosen the rule acted in accordance with predictions 3 and 4, they will continue to request '4' thereafter.*

3.2 Alternative predictions

Let us now ask what kind of results can be expected from application of models based on a different characterisation of players' motivational complexity. Before we discuss the experimental results, we will answer this question by briefly considering two alternative

approaches: the “inequity aversion” theory devised by Fehr and Schmidt (1999) and the model of reciprocity introduced by Rabin (1993).

3.2.1 *Inequity aversion in the Exclusion Game*

The idea at the basis of Fehr and Schmidt’s (1999) model is that agents suffer from both advantageous and (to a greater extent) disadvantageous inequality. Mathematically, given n players, the utility function of player i is:

$$u_i(x_1, x_2, \dots, x_n) = x_i - \frac{\alpha_i}{n-1} \sum_{j \neq i} \max\{x_j - x_i, 0\} - \frac{\beta_i}{n-1} \sum_{j \neq i} \max\{x_i - x_j, 0\}$$

where α_i and β_i are the weights of utility losses from disadvantageous (α_i) and advantageous (β_i) inequality, with $\beta_i \leq \alpha_i$ and $0 \leq \beta_i < 1$.

Predictions about the solution of the Exclusion Game depend on the players with whom each active player compares his/her payoff. If s/he cares only about the other active player, the unique Nash equilibrium of the game is the one in which both the active players request €6. But, if the third player’s payoff enters the utility function of the active players, then (4,4) is a Nash equilibrium if $\beta_1 = \beta_2 > 2/3$.⁹

However, the distinguishing feature of this model is that according to it there should be no differences at all between the way the game is played in the first phase and the way it is played in the third phase of the experiment. In other words, phase 2 is irrelevant to this model.

3.2.2 *Reciprocity in the Exclusion Game*

Rabin (1993) conceives a model of choice in which the agent’s utility depends both on his material payoff (π_i) and on a measure of reciprocal kindness (Rabin, 1993 :1286-1287). In particular, player i judges his own kindness f_i toward j , given his strategy a_i and his beliefs b_j about j ’s strategy, as the distance between the payoff he gives to j and an equitable payoff – which is defined as the average between the maximum and the minimum payoffs that he

⁹ The strategy vector (4,4) is a Nash equilibrium if

$$u_1(4,4,4) > u_1(6,4,2) \Leftrightarrow 4 > 6 - 3\beta_1$$

and

$$u_2(4,4,4) > u_2(4,6,2) \Leftrightarrow 4 > 6 - 3\beta_2$$

could give to j . Kindness of j toward i (defined as f_i'), given i 's beliefs about j 's beliefs about i 's strategy (c_i), is measured as the difference between the payoff that j gives to i and the equitable payoff (calculated like before).

The utility function of player i is given by

$$U_i(a_i, b_j, c_i) \equiv \pi_i(a_i, b_j) + f_j'(b_j, c_i) [1 + f_i'(a_i, b_j)]$$

Note that if i perceives the action by j as not kind, he will gain utility by reacting with an unkind action, while if he perceives j 's choice as kind, he will gain utility by reciprocating with a kind action. Given these preferences, and using the tools of psychological game theory, Rabin introduces the concept of *fairness equilibrium*, which is defined by a vector of strategies which are mutually best responses to each other and by a set of beliefs compatible with those strategies (Rabin; 1993: 1288).

Rabin provides a theoretical framework that has also proved very useful in developing the model of choice subject to the experimentation described here (Grimalda and Sacconi 2002, 2005). To be noted, however, is that the basic difference between the notions of fairness incorporated into the indexes of 'conformity' (Sacconi and Grimalda) or 'kindness' (Rabin) used by the two theories respectively generates a sharp difference in the players' expected behaviours in our experimental game. Assuming Rabin's model of motivation we would predict (6,6) as the unique fairness equilibrium of the Exclusion Game, both in phase one and three. To see this, consider the fact that in the Exclusion Games active players are the dictators in a standard Dictator Game. Each of them is endowed with a sum of R/2 and has to decide how much of that sum to give to the third player; but his own payoff does not depend upon the action of the other active player. As a consequence, the two kindness measures are both equal to zero.

This is due to the fact that what Rabin understands to be fairness is in fact a direct and reciprocal attitude of kindness shown by the actions that players undertake directly toward each other and vice versa. This is a personal relationship and a personally orientated attitude. By contrast, the model of conformist preferences is based on an impersonal and impartial principle of fairness chosen under an hypothesis of symmetry and exchangeability of personal places, which must be applied to whichever player is affected by the real game in whatever

role only because he may have been party to the constitutional of choice of rules. In so far as impersonality and impartiality of the principle (rule) are distinctive of this approach, it can be called a theory of 'ideal preferences' based on 'ethical principles of fairness'. For the same reason, it is doubtful that 'fairness' (a typical ethical concept) is the true feature of Rabin's model - which can be more properly understood as a 'kindness' model (where kindness is a feature of direct personal relationships)

The same prediction – both active players will choose strategy '6' - would derive from the application of any other model where 'intentionality' is construed as an *attitude* (in the player's action) directed toward another player, who then perceives the intention of the opponent's action as being to generate a favourable or unfavourable consequence specifically to himself, like a manifestation of kindness, friendship or hostility - that is, any model based on the idea that a player's choice is influenced by the perceived intentions of the opponent directed toward the player himself and by the willingness to react to these intentions - see for example Falk and Fischbacher, 2001.¹⁰)

3.3 Results¹¹

3.3.1 General description of behaviours

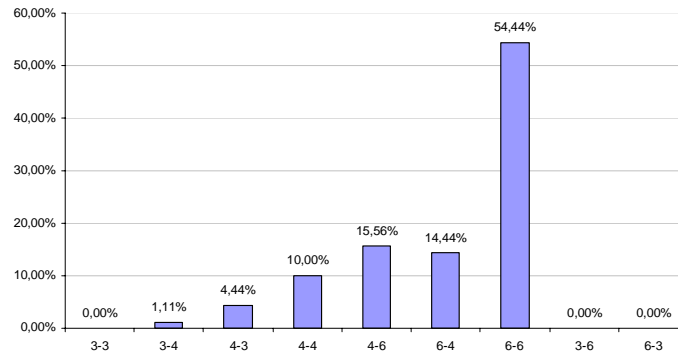
We begin with a general description of the subjects' behaviour. In order to be clear, let us remind that each subject plays two *rounds* of the game in phase 1, and then (if he is not selected to be a dummy player) he plays again *twice* (we call these '*times*') the basic game in phase three. In the first phase subjects played once in the G1 role and once in the G2 role. Considering all the 90 subjects involved in the three sessions (Figure 7¹²), we observed that a large majority of players (54.4%) chose to ask for six euros when playing in both the G1 and G2 roles, leaving nothing for the third player.

¹⁰ It should be noted that this is not the only way intentionality can be construed: typically we say that deontological reasoning is based on intention and not on consequences, because it is driven by the 'nature' (rightness or wrongness) of the action itself, without any reference to its consequences for anybody. Rightness or wrongness are typically decided in terms of the consistency of actions with some abstract moral principle. In this case intentionality is expressed by the desire to act according with an impersonal and impartial principle or - to say it differently - the intention to act rightly, without any consideration of the consequence resulting to any potential player whatsoever.

¹¹ See Appendix 3 for the complete dataset.

¹² The two numbers represent the choices made in the G1 and G2 roles respectively.

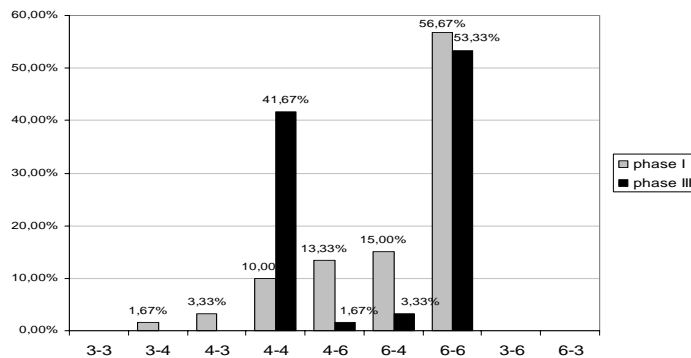
Figure 7: Players' choice in phase 1 (n=90)



About one third of the subjects made different choices in the two symmetrical roles G1 and G2, choosing to ask for 6 euros as G1 (G2) and 4 euros as G2 (G1). Although we cannot completely exclude the possibility of confusion, a look at the questionnaires suggests that this particular behaviour can be explained mostly in terms of regret, which induced subjects to adjust, in the second round, the choice made in the first round.

Jumping for a moment to phase three, it will be seen that some players were assigned the role of the dummy player. As a consequence, because they never played during this phase, their behavior cannot be compared with their behaviour in the first phase. The remaining players instead took an active role in both rounds of the third phase. Hence, if we want to compare the results of phase one and phase three, we must limit our analysis to the 60 subjects that played in G1 or G2 in this phase (Figure 8).¹³

Figure 8: Active players' choices in phase I and phase III (n=60)



¹³ With regard to phase three, two numbers represent the choices at time 1 and at time 2 respectively.

Three main facts emerge from comparison between the two distributions: first, moving from phase 1 to phase 3 there is a significant growth of the share of subjects who asked for 4 euros at both the first and the second time they played in phase 3. Second, the percentage of subjects that chose to keep the maximum amount for themselves (6 euros) in both the rounds of the first phase did not greatly differ from the percentage of players who decided to ask for 6 euros in both the two rounds of the third phase. Finally, we can observe a shift from the unimodal distribution of the first phase's choice to the clearly bimodal distribution of the choices of the third phase, in which only a negligible number of subjects asked for different amounts in the two rounds. Table 1 gives a more detailed picture of this polarization between the "4,4" and the "6,6" choices.¹⁴

Table 1. Choices in phase I and III (N=60)

		Phase III							
		3-3	3-4	4-3	4-4	4-6	6-4	6-6	
Phase I	3-3	0	0	0	0	0	0	0	0
	3-4	0	0	0	1	0	0	0	1
	4-3	0	0	0	1	0	0	1	2
	4-4	0	0	0	4	0	0	2	6
	4-6	0	0	0	3	0	0	5	8
	6-4	0	0	0	4	0	0	5	9
	6-6	0	0	0	12	1	1	20	34
		0	0	0	25	1	1	33	60

¹⁴ To test whether a difference exists between the choices in phase I (before agreement on the rule) and those in phase three (after agreement on the rule) we can distinguish between equity oriented choices, corresponding to "Class E= ask for '4' at least once, but do not ask for '6' " and more self-interest-oriented choices, corresponding to "Class S= asks for '6' at least once". We can introduce the null hypothesis that the number of subjects that move from Class E to Class S is the same as the number of subjects moving from Class S to Class E.

		PHASE III		
		E	S	
PHASE I	E	6	3	9
	S	19	32	51
		25	35	60

On looking at the joint distribution of frequency, we can reject the null hypothesis that the probabilities of being in cells [S,E] and [E,S] are the same (McNemar's Chi-squared =10,22, df=1, p-value=0,0013).

3.3.2 Choice and implementation of the division rule

Inspection of the data in phase two, when the subjects are requested to agree over a principle and a division rule, shows that 19 of the 30 groups (57 subjects) chose principle 1 and rule (33%, 33%, 33%), 10 groups (30 subjects) chose principle 2 and rule (50%, 50%, 0%) and only one group (3 subjects) chose principle 2 and rule (50%, 33%, 17%).

Unanimity on a general principle for each group was reached within a maximum of three trials, with a large majority of groups reaching agreement in the first trial. The maximum number of trials required to reach unanimous agreement on a specific division rule was seven,¹⁵ but most groups did not go beyond the first trial.

With regard to active players, 24 of the 38 who chose the (33%, 33%, 33%) rule decided to implement it at both times they played in the third phase. Rule (50%, 50%, 0%) was implemented by 19 of the 20 active players belonging to the groups that selected that rule (Table 2). Only one of the two members with an active role, who agreed on the (50%, 33%, 17%) rule, decided to implement it at both times of third phase (Table 2).

Table 2. Division rule choice and implementation (Active players; N=30)

<i>RULE</i>	<i>N.</i>	<i>Implementation of the rule at first time (freq.)</i>	<i>Belief about other's conformity (time 1) (freq.)</i>	<i>Implementation of the rule at second time (freq.)</i>	<i>Beliefs about other's conformity (time 2) (freq.)</i>	<i>Implementation of the rule in both time 1 and 2 (freq.)</i>
{33%, 33%, 33%}	38	25	23	25	24	24
{50%, 50%, 0%}	20	19	19	19	19	19
{50%, 33%, 17%}	2	1	1	1	0	1
<i>TOT.</i>	60	45	43	45	43	44

At both times of phase three, almost all the players that implemented the rule (both the rules) predicted an outcome of the game compatible with reciprocal conformity. On the other hand,

¹⁵ Reached only by the members of the group that agreed on the (50%, 33%, 17%) rule.

almost all the players who decided not to conform with the rule predicted that their opponents would do the same (see Appendix 3 for details).

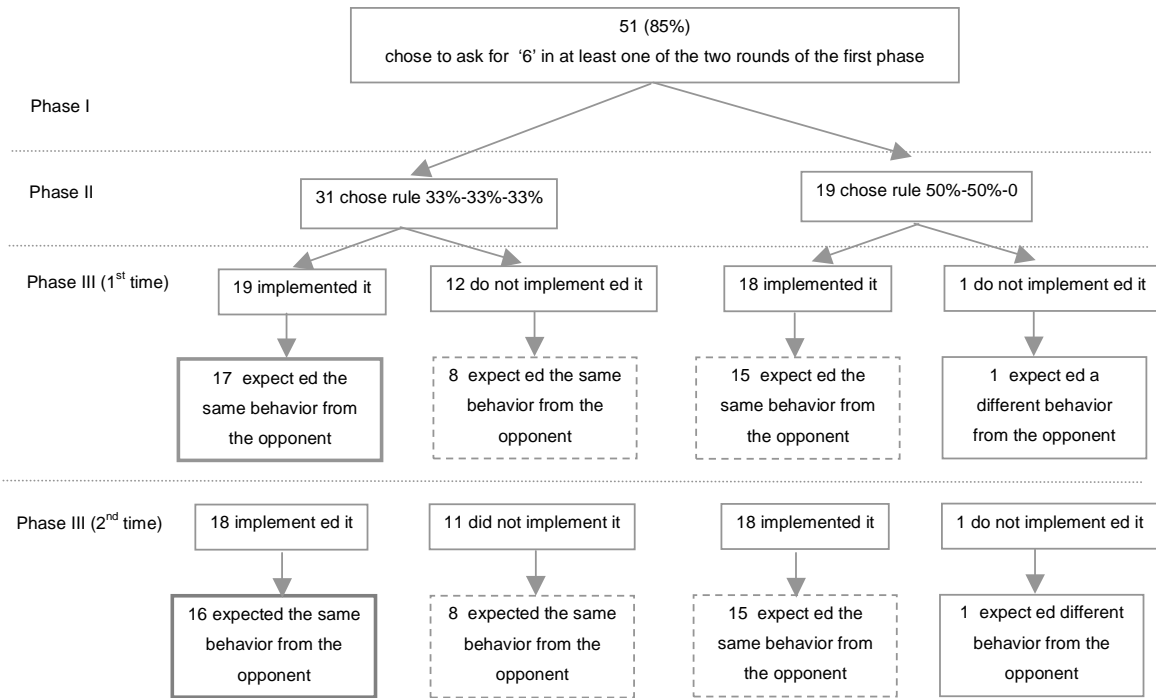
Forty four players (73%) made correct predictions about the willingness of the opponent to conform with the rule. As shown by the data, information about the percentage of players that chose to implement the rule the first time, had no effect on their choice at the second time of this phase. It is not troubling however, in that this information largely confirmed their initial prediction (those who predicted conformity first time conformed and then received the information that most of other players had conformed), so that there was no reason to make a change in their beliefs.

4. Identifying the players with conformist preferences

The most interesting fact emerging from these data is the significant *change* from a situation in which most players chose a strategy compatible with the pursuit of their self interest, asking for half of the sum for themselves and leaving nothing for the third player, to a situation in which there was a significant growth in the propensity to choose an equal division by implementing a widely shared rule of equal distribution of the benefits. This shift becomes particularly clear if, considering only the 60 subjects with an active role in the third phase, we look at the behaviour of the 51 of them (85%) that asked for 6 euros in at least one of the rounds of the first phase, so taking a decision that cannot discriminate between the two self-interest and conformistic preferences hypotheses (Figure 9). Indeed, we can observe that 31 of them chose the (33%, 33%, 33%) rule in the second phase, 19 of these implemented it in the third phase, with 17 believing that their opponents would do the same. Eight of the remaining 11 subjects who decided to ask for the maximum amount did not believe that their opponents would implement the rule.

17 of the 18 players who agreed on the (50%, 50%, 0) rule implemented it expecting the same behaviour by their opponents.

Figure 9: The choices of players who chose to ask for '6' at least once in the first phase



Following the choices of this class of subjects through the three phases of the experiment, we can conclude that at least 16 of them (one third of all the active players who asked for 6 euros in at least one of the rounds of the first phase) made choices that can be considered perfectly compatible with our hypotheses (boxes with thick borders in figure 9). Indeed, the change in their behaviour from phase 1 to phase 2 can be attributed to the activation of conformist motives due to the introduction of a shared non-binding rule of division and to the emergence of beliefs about reciprocal conformity.

In addition, 8 of the 12 (11 in the second round) players that decided not to implement the 33%, 33%, 33% rule did so believing that the counterpart would do the same.

This last observation does not conflict with the conformist preferences hypothesis. However, because it is also compatible with the assumption that players are characterized by self-interest, it does not allow us to identify pure-conformist players. The same can be said with regard to the 15 players who decided to implement the (50%, 50%, 0) rule expecting the same behaviour by their opponents (boxes with dotted borders in figure 9).

5. Conclusions

We can summarize our results by saying that at least one third of the players who were always active behaved and asked for 6 euros in at least one of the rounds of the first phase were motivated by conformist preferences, while the remaining subjects can be identified as being either self-interested or conformist. Only four players displayed behavior compatible with the inequity-averse hypothesis. In particular, with respect to our empirical predictions:

- i) The observation that most players chose strategy '6' in both the rounds of the first phase is compatible with prediction 1.
- ii) The fact that, in the second phase, a large number of players agreed on the (33%, 33%, 33%) rule is compatible with predictions 2.
- iii) For a significant number of subjects, having agreed on a rule seems to have been sufficient reason to generate expectations about reciprocal conformity.¹⁶
- iv) There is a close correlation between a player's belief about the opponent's willingness to conform with the rule and his decision to implement it (prediction 3).
- v) A *significant* number of those players who egoistically chose strategy '6' in at least one of the rounds of the first phase and who agreed on the rule of equal division in phase two, decided to implement the rule in phase three; and these are definitely *most* of those who, having acted as just described, entertained the belief that the agreed rule would have been played by their counterpart; all this is in accordance with prediction 4.
- vi) Information that confirms beliefs about conformity does not change the willingness to conform, in accordance with prediction 5.

These results allow us concluding that the experimental evidence backs the Sacconi and Grimalda's model of choice based on conformist preferences and comprehensive utility functions. The significant shift in the behaviour of the players in the transition from phase one to phase three is strictly consistent with the hypothesis that, having these subject realised that the constitutional nature of the choice in phase two asked for a fairness principle of conduct incompatible with their behaviour in phase 1, the very fact of having agreed on that principle

¹⁶ This hypothesis is supported by the replies to the questionnaire.

activated their conformist motivation to conform with it, granted that they believed that the same principle was conformed by the counterparties. Notice that, whereas the players changing their behaviour crucially corroborate the theory, also most of players that exhibit the same behaviour in phase one and three (given their choice in phase two and their beliefs), are consistent with the theory (i.e. do not provide any anomaly to the model). This cannot be said for the concurrent models we considered, which cannot explain our data. With regard to these alternative theories, in fact, we can conclude that:

a) Models of inequity-aversion fail to explain the observation of different behaviours in phase one and in phase three by the same subjects. Why should the subjects be inequity averse in the third phase if they were not so motivated in the first? This finding may be explained within the “inequity –aversion” framework by saying that the introduction of phase two, modelled as a constitutional choice, induced a change in the definition of the reference group, inducing subjects that in the first phase did not consider the payoff of the dummy as relevant to include it instead in the third phase. However, this explanation would make the inequity-aversion approach closely akin to the conformist preference model, where players that consider the advantage connected to the constitution of a social union (that is, the advantage of being allowed to play the game because of agreement on a rule) develop the motivational basis for a change in their behaviour (but consider that we found that players’ action also depended on their beliefs concerning the reciprocity of the counterpart, an aspect that does not have any significant role in the inequity-aversion model).

b) ‘Direct reciprocity’ models, or reciprocity models based on ‘direct kindness’, fail to predict the dramatic change in the behaviour pattern shown by subjects between the first phase, when it in fact accords with the direct kindness prediction, to the third one, when it diverges substantially. Note that dummy players did not change their status during the process that took the player from phase one to phase three, i.e. at the last step there were still dummy players, and they could not manifest direct attitude or intention toward the active players since they simply do not make decisions. Nevertheless, having a rule in mind, one that has been agreed even if not binding or exogenously enforced, seems to be enough to change the players’ behaviour significantly. This suggests that some sort of commitment to the principle itself, and beliefs concerning reciprocal conformity with it, has motivational effects. Quite paradoxically, in the situation under experimentation ‘fairness’, understood according to Rabin model as direct reciprocity between the two active players, would imply to discriminate

against the weak player, and would result in a behaviour completely indistinguishable from the conduct that shares all the pie amongst the strong players. Players who are fair according to our model act against 'fairness' in Rabin's sense, and should elicit a punitive response by the other payers if Rabin's model were true. On the contrary our explanation accounts for these behaviours in terms of a desire - grounded on what can be seen as a reciprocal 'sense of justice' - for conformity to an impartial principle of fairness agreed under a veil of ignorance, what is much more coherent to the notion of justice given in normative economics and political philosophy (Rawls, 1971).

References:

- Bolton, G.E. (1991), "A Comparative Model of Bargaining: Theory and Evidence" *American Economic Review*, 81, 1096-1136.
- Bolton, G.E. and Ockenfels, A. (2000), "A theory of Equity, Reciprocity and Competition". *American Economic Review*, 100, 166-193.
- Cubitt, R., C. Starmer, and R. Sugden (1998), "On the Validity of the Random Lottery Incentive system", *Experimental Economics*, 1: 115:131.
- Falk, Armin and Fischbacher, Urs (2001), "A Theory of Reciprocity" *Institute for Empirical Research in Economics Working Paper No. 6*.
- Fehr, E, and Schmidt, K. M. (1999). "A Theory of Fairness, Competition and Co-operation." *Quarterly Journal of Economics*, 114, 817-868.
- Geanakoplos, J. Pearce, D. and Stacchetti, E. (1989), "Psychological Games and Sequential Rationality", *Games and Economic Behavior*, Vol. 1, pp. 60-79.
- Grimalda, G. and L. Sacconi, (2002), "The constitution of the non profit enterprise, ideals, conformism and reciprocity", *University Carlo Cattaneo - LIUC paper n.155*.
- Grimalda, G., L. Sacconi (2005a) "The Constitution of the Not-for-Profit Organisation: Reciprocal Conformity to Morality" Forthcoming in *Constitutional Political Economy*, Vol 16(3), September 2005.
- Sacconi L. and Grimalda G. (2005b) "Ideals, conformism and reciprocity: A model of Individual Choice with Conformist Motivations, and an Application to the Not-for-Profit Case" in (L.Bruni and P.L.Porta eds.) *Handbook of Happiness in Economics*, Edward Elgar, London, (forthcoming)
- Rabin M. (1993), "Incorporating Fairness into Game Theory and Economics" *The American Economic Review*, 83(5):1281-1302.
- Rawls J. (1971), *A Theory of Justice*, Oxford U.P, Oxford.
- Starmer, C., and R. Sugden (1991) "Does the Random-Lottery Incentive System Elicit True Preferences? An Experimental Investigation", *American Economic Review*, 81: 971-978.

Appendix 1. The Psychological Nash Equilibrium ⁺⁺

At the basis of this concept of solution is the idea that, in equilibrium, rational players' beliefs must be coherent with their strategies. For example, if in equilibrium player i observes that j plays the strategy $\sigma_j \in \Sigma_j$, then i 's 1st order belief must assign probability 1 to the fact that j plays that strategy with probability 1, and 0 to the other strategies. Furthermore, given the equilibrium strategy σ_i of player i , and given the usual assumption of common knowledge of players' rationality, i must expect j to have the same rational beliefs, coherently with the fact that i plays σ_i with probability 1. This means that i 's 2nd order beliefs must assign probability 1 to the fact that j believes that i plays σ_i . More generally, in equilibrium, all the 1st order beliefs must be single-point distributions assigning probability 1 to the equilibrium strategy. The higher order beliefs must be consistent with this condition and with the assumption of common knowledge of players' rationality (Geneakoplos et al., 1989:64). Let us call $\beta_i(\sigma)$ the distribution of beliefs that satisfy this coherence condition over σ and $\beta(\sigma) = (\beta_1(\sigma), \dots, \beta_n(\sigma)) \in \beta$ the profile of such beliefs for the n players.

Recalling the definition of b_i as the vector of beliefs of each order for player i , and of $b = (b_1, \dots, b_n)$ as the profile of beliefs for each of the n players, we can define the Psychological Nash Equilibrium (Geneakoplos et al., 1989:65):

A Psychological Nash Equilibrium for n -player normal form game is a pair $(\hat{b}, \hat{\sigma}) \in \beta \times \Sigma$ such that:

i) $\hat{b} = \beta(\hat{\sigma})$

ii) for each $i \in I$ and $\sigma_i \in \Sigma_i$, $V_i(\hat{b}_i, (\sigma_i, \sigma_{-i})) \leq V_i(\hat{b}_i, \hat{\sigma}_i)$.

Condition (ii) is a restatement of the standard Nash equilibrium condition affirming that the equilibrium strategy must assign a payoff no smaller than the ones attained by any other feasible strategy, given the opponent's strategy and the beliefs. Condition (i) requires the beliefs to be coherent with the equilibrium strategy. Note that if beliefs are not part of the

⁺⁺ This appendix draws on Sacconi and Grimalda (2003)

utility function then condition (i) becomes redundant and the definition boils down to the standard Nash equilibrium definition.

Appendix 2. Voting Procedure

Subjects were randomly assigned to five groups with three members (identified with the numbers from 1 to 5). Each member could read the number of his group on his computer screen but could not interact with the other members, nor identify them.

The experimenter distributed a form like the one in figure 1a

Figure 1a: Form for the general principle selection

EXPERIMENT PREFCOMP (PILOT) 14/10/2004

**SECOND PHASE
GENERAL PRINCIPLE CHOICE**

PRINCIPLE 1

Every player should share benefits, independently from his role. In particular, who has not the possibility to choose should not restrict the others choices.

PRINCIPLE 2

People who play under a dominant role could obtain a higher share of benefits.

PLAYER ID:..... GROUP NUMBER:.....

N	Player's own choice	Other players (please, do not fill)			
1.	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	AGREE NOT AGREE
2.	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	AGREE NOT AGREE
3.	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	AGREE NOT AGREE
4.	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	AGREE NOT AGREE
5.	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	<input type="checkbox"/> 1 <input type="checkbox"/> 2	AGREE NOT AGREE

OUTCOME: 1 2

The subjects were asked to fill in the “ ID ” and “Group’s number” box and to select their preferred principle by ticking one of the two boxes in the “Player’s choice” column. The experimenters collected the forms and checked the votes, writing on each player’s form the choices of the other members of the group. If the members of some groups did not reach unanimous agreement, the experimenter again distributed the forms to all the subjects.¹⁷

¹⁷ This made it impossible to identify the members of a particular group exploiting the information about the outcome of the voting procedure.

Members of the groups that did not reach agreement were asked to vote again, while the others had to wait. The experimenter collected the forms, checked the votes and repeated the same procedure until all the groups had reached agreement. The maximum number of trials allowed was five.

After the voting for selection of the general principle new forms like the ones in figures 2a and 3a were distributed. These stated particular division rules deduced from the general principle. Each subject received a form stating the rules deduced from the principle selected in the previous stage. The voting procedure was the same as the one adopted for the principles selections, but the maximum number of trials was now ten.

Figure 2a: Form for the selection of rule deduced from principle 1
 Figure 3a: Form for the selection of rule deduced from principle 2

EXPERIMENT PREFCONF (PILOT) 14/10/2004

SECONDA PHASE DIVISION RULE SELECTION

PRINCIPLE 1
 "Every player should share benefits, independently from his role. In particular, who has not the possibility to choose should not receive less than others."

N	O1 (ACTIVE)	O2 (ACTIVE)	O3 (NOT ACTIVE)
1	33% (4)	25% (3)	42% (5)
2	33% (4)	33% (4)	33% (4)
3	33% (3)	33% (4)	40% (5)

PLAYER'S ID: GROUP NUMBER:

N	Player's own choice	Other players (please do not fill)		
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				

OUTCOME

EXPERIMENT PREFCONF (PILOT) 14/10/2004

SECONDA PHASE DIVISION RULE SELECTION

PRINCIPLE 2
 "People who play under a decisional role could obtain a higher share of benefits."

N	O1 (ACTIVE)	O2 (ACTIVE)	O3 (NOT ACTIVE)
1	30% (3)	33% (4)	37% (5)
2	33% (4)	30% (3)	37% (5)
3	30% (3)	30% (3)	40% (5)

PLAYER'S ID: GROUP NUMBER:

N	Player's own choice	Other players (please do not fill)		
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				

OUTCOME

At the end of the voting procedure, the experimenter inserted the rule selected in a form that appeared on the screen of each subject.

Appendix 4. The questionnaire.

1) What strategy did you choose in phase one?

As G1..... ; As G2.....

Could you describe the reasons for these choices?

2) On what rule did you reach agreement in the second phase? To what extent do you agree with? this rule?

3) What were your expectations about the willingness of the other player to apply the rule in the third phase (Or the G1 and G2 willingness if your role was G3) ?

4) What have been the main differences between the first and the third phase?

Elenco dei papers del Dipartimento di Economia

- 2000.1 *A two-sector model of the effects of wage compression on unemployment and industry distribution of employment*, by Luigi Bonatti
- 2000.2 *From Kuwait to Kosovo: What have we learned? Reflections on globalization and peace*, by Roberto Tamborini
- 2000.3 *Metodo e valutazione in economia. Dall'apriorismo a Friedman*, by Matteo Motterlini
- 2000.4 *Under tertiarisation and unemployment*. by Maurizio Pugno
- 2001.1 *Growth and Monetary Rules in a Model with Competitive Labor Markets*, by Luigi Bonatti.
- 2001.2 *Profit Versus Non-Profit Firms in the Service Sector: an Analysis of the Employment and Welfare Implications*, by Luigi Bonatti, Carlo Borzaga and Luigi Mittone.
- 2001.3 *Statistical Economic Approach to Mixed Stock-Flows Dynamic Models in Macroeconomics*, by Bernardo Maggi and Giuseppe Espa.
- 2001.4 *The monetary transmission mechanism in Italy: The credit channel and a missing ring*, by Riccardo Fiorentini and Roberto Tamborini.
- 2001.5 *Vat evasion: an experimental approach*, by Luigi Mittone
- 2001.6 *Decomposability and Modularity of Economic Interactions*, by Luigi Marengo, Corrado Pasquali and Marco Valente.
- 2001.7 *Unbalanced Growth and Women's Homework*, by Maurizio Pugno
- 2002.1 *The Underground Economy and the Underdevelopment Trap*, by Maria Rosaria Carillo and Maurizio Pugno.
- 2002.2 *Interregional Income Redistribution and Convergence in a Model with Perfect Capital Mobility and Unionized Labor Markets*, by Luigi Bonatti.
- 2002.3 *Firms' bankruptcy and turnover in a macroeconomy*, by Marco Bee, Giuseppe Espa and Roberto Tamborini.
- 2002.4 *One "monetary giant" with many "fiscal dwarfs": the efficiency of macroeconomic stabilization policies in the European Monetary Union*, by Roberto Tamborini.
- 2002.5 *The Boom that never was? Latin American Loans in London 1822-1825*, by Giorgio Fodor.

2002.6 *L'economia senza banditore di Axel Leijonhufvud: le 'forze oscure del tempo e dell'ignoranza' e la complessità del coordinamento*, by Elisabetta De Antoni.

2002.7 *Why is Trade between the European Union and the Transition Economies Vertical?*, by Hubert Gabrisch and Maria Luigia Segnana.

2003.1 *The service paradox and endogenous economic growth*, by Maurizio Pugno.

2003.2 *Mappe di probabilità di sito archeologico: un passo avanti*, di Giuseppe Espa, Roberto Benedetti, Anna De Meo e Salvatore Espa.
(*Probability maps of archaeological site location: one step beyond*, by Giuseppe Espa, Roberto Benedetti, Anna De Meo and Salvatore Espa).

2003.3 *The Long Swings in Economic Understanding*, by Axel Leijonhufvud.

2003.4 *Dinamica strutturale e occupazione nei servizi*, di Giulia Felice.

2003.5 *The Desirable Organizational Structure for Evolutionary Firms in Static Landscapes*, by Nicolás Garrido.

2003.6 *The Financial Markets and Wealth Effects on Consumption An Experimental Analysis*, by Matteo Ploner.

2003.7 *Essays on Computable Economics, Methodology and the Philosophy of Science*, by Kumaraswamy Velupillai.

2003.8 *Economics and the Complexity Vision: Chimerical Partners or Elysian Adventurers?*, by Kumaraswamy Velupillai.

2003.9 *Contratto d'area cooperativo contro il rischio sistemico di produzione in agricoltura*, di Luciano Pilati e Vasco Boatto.

2003.10 *Il contratto della docenza universitaria. Un problema multi-tasking*, di Roberto Tamborini.

2004.1 *Razionalità e motivazioni affettive: nuove idee dalla neurobiologia e psichiatria per la teoria economica?* di Maurizio Pugno.
(*Rationality and affective motivations: new ideas from neurobiology and psychiatry for economic theory?* by Maurizio Pugno.

2004.2 *The economic consequences of Mr. G. W. Bush's foreign policy. Can th US afford it?* by Roberto Tamborini

2004.3 *Fighting Poverty as a Worldwide Goal* by Rubens Ricuperò

- 2004.4 *Commodity Prices and Debt Sustainability* by Christopher L. Gilbert and Alexandra Tabova
- 2004.5 *A Primer on the Tools and Concepts of Computable Economics* by K. Vela Velupillai
- 2004.6 *The Unreasonable Ineffectiveness of Mathematics in Economics* by Vela K. Velupillai
- 2004.7 *Hicksian Visions and Vignettes on (Non-Linear) Trade Cycle Theories* by Vela K. Velupillai
- 2004.8 *Trade, inequality and pro-poor growth: Two perspectives, one message?* By Gabriella Berloff and Maria Luigia Segnana
- 2004.9 *Worker involvement in entrepreneurial nonprofit organizations. Toward a new assessment of workers? Perceived satisfaction and fairness* by Carlo Borzaga and Ermanno Tortia.
- 2004.10 *A Social Contract Account for CSR as Extended Model of Corporate Governance (Part I): Rational Bargaining and Justification* by Lorenzo Sacconi
- 2004.11 *A Social Contract Account for CSR as Extended Model of Corporate Governance (Part II): Compliance, Reputation and Reciprocity* by Lorenzo Sacconi
- 2004.12 *A Fuzzy Logic and Default Reasoning Model of Social Norm and Equilibrium Selection in Games under Unforeseen Contingencies* by Lorenzo Sacconi and Stefano Moretti
- 2004.13 *The Constitution of the Not-For-Profit Organisation: Reciprocal Conformity to Morality* by Gianluca Grimalda and Lorenzo Sacconi
- 2005.1 *The happiness paradox: a formal explanation from psycho-economics* by Maurizio Pugno
- 2005.2 *Euro Bonds: in Search of Financial Spillovers* by Stefano Schiavo
- 2005.3 *On Maximum Likelihood Estimation of Operational Loss Distributions* by Marco Bee
- 2005.4 *An enclave-led model growth: the structural problem of informality persistence in Latin America* by Mario Cimoli, Annalisa Primi and Maurizio Pugno
- 2005.5 *A tree-based approach to forming strata in multipurpose business surveys*, Roberto Benedetti, Giuseppe Espa and Giovanni Lafratta.

2005.6 *Price Discovery in the Aluminium Market* by Isabel Figuerola-Ferretti and Christopher L. Gilbert.

2005.7 *How is Futures Trading Affected by the Move to a Computerized Trading System? Lessons from the LIFFE FTSE 100 Contract* by Christopher L. Gilbert and Herbert A. Rijken.

2005.8 *Can We Link Concessional Debt Service to Commodity Prices?* By Christopher L. Gilbert and Alexandra Tabova

2005.9 *On the feasibility and desirability of GDP-indexed concessional lending* by Alexandra Tabova.

2005.10 *Un modello finanziario di breve periodo per il settore statale italiano: l'analisi relativa al contesto pre-unione monetaria* by Bernardo Maggi e Giuseppe Espa.

2005.11 *Why does money matter? A structural analysis of monetary policy, credit and aggregate supply effects in Italy*, Giuliana Passamani and Roberto Tamborini

2005.12 *Conformity and Reciprocity in the "Exclusion Game": an Experimental Investigation* by Lorenzo Sacconi and Marco Faillo

PUBBLICAZIONE REGISTRATA PRESSO IL TRIBUNALE DI TRENTO