

MPRA

Munich Personal RePEc Archive

Ashamed to be Selfish

Dillenberger, David and Sadowski, Philipp

16 April 2008

Online at <https://mpra.ub.uni-muenchen.de/8343/>
MPRA Paper No. 8343, posted 20 Apr 2008 05:34 UTC

Ashamed to be Selfish*

David Dillenberger and Philipp Sadowski[†]

April 17, 2008

Abstract

We study a two-stage choice problem, where alternatives are allocations between the decision maker (DM) and a passive recipient. The recipient observes choice behavior in stage two, while stage one choice is unobserved. Choosing selfishly in stage two, in the face of a fairer available alternative, may inflict shame on DM. DM has preferences over sets of alternatives that represent period two choices. We axiomatize a representation that identifies DM's selfish ranking, her norm of fairness and shame. Altruism is the most prominent motive that can explain non-selfish choice. We identify a condition under which shame to be selfish can mimic altruism, when only stage-two choice is observed by the experimenter. An additional condition implies that the norm of fairness can be characterized as the Nash solution of a bargaining game induced by the second-stage choice problem. The representation is generalized to allow for finitely many recipients and applied to explain a social decision maker's incentive for obfuscation.

JEL Classifications: C78, D63, D64, D78, D80, D81

1. Introduction

1.1. Motivation

The notions of fairness and altruism have attracted the attention of economists in different contexts. The relevance of these motives to decision making is both intuitively convincing and well documented. For example in a classic “dictator game,” where one person gets to anonymously divide, say, \$10 between herself and a partner, people tend not to take the whole amount for themselves, but to give a sum of between \$0 and \$5 to the other player. They act as if they are trading off a concern for fairness or for the other person's incremental

*Preliminary. We thank Roland Benabou and Wolfgang Pesendorfer for their invaluable support. We are also grateful to Eric Maskin, Stephen Morris, Charles Roddie and Tymon Tatur for helpful suggestions. This paper was written in part while the authors were graduate fellows at the University Center for Human Values, Princeton University.

[†]Department of Economics, Princeton University, Princeton, NJ 08544,
ddillenb@princeton.edu and sadowski@princeton.edu.

wealth and a concern for their own.¹ Thus, preferences for fairness as well as preferences for altruism have been suggested and considered (for example Fehr and Schmidt [1999], Anderoni and Miller [2002], and Charness and Rabin [2002]).

Recent experiments, however, show that this interpretation may be rash: Dana, Cain and Dawes (2006) study a variant of the same dictator game, where the dictator is given the option to exit the game before the recipient learns it is being played. If she opts out, she is given a specified amount of money and the recipient gets nothing, as the game has not taken place. It turns out that about a third of the participants choose to leave the game when offered \$9 for themselves and \$0 for the recipient. Write this allocation as (\$9, \$0). Such behavior contradicts altruistic concern regarding the recipient's payoff, because then the allocation (\$9, \$1) should be strictly preferred. It also contradicts purely selfish preferences, as (\$10, \$0) would be preferred to (\$9, \$0). Instead, people seem to suffer from behaving egoistically in a choice situation where they could dictate a fairer allocation. Hence, if they can avoid getting into such a situation, they happily do so. Real-life scenarios with this character could be:

- donating to a charity over the phone but wishing not to have been home when the call came,
- crossing the road to avoid meeting a beggar.

Our explanation of this type of behavior is the following: Whether a person's actions are observed or not plays a crucial role in determining her behavior. We term "shame" the motive that distinguishes choice behavior when observed from choice behavior when not observed. In our model, individuals are selfish when not observed. Thus, concern for another person's payoff is motivated not by altruism, but by avoiding the feeling of shame that comes from behaving selfishly when observed.² The interpretation is that, if people are observed, they feel shame when they do not choose the fairest available alternative.³

We axiomatically formalize the notion of shame and its interaction with selfishness as described above. To this end, we consider games like the one conceived by Dana et al (2006) as a two-stage choice problem. In the first stage, the decision maker (DM) chooses a "menu," a set of payoff-allocations between herself and the anonymous recipient. This choice is not observed by the recipient. In the second stage, she makes a potentially anonymous choice from the alternatives on this menu, where the recipient observes the chosen alternative in full

¹See for example Camerer (2003).

²To distinguish shame from guilt, note that guilt is typically understood to involve regret, even in private, while, according to Buss (1980), "*shame is essentially public; if no one else knows, there is no basis for shame. [...] Thus, shame does not lead to self-control in private.*" We adopt the interpretation that even observation of a selfish behavior without identification of its purveyor can cause shame.

³In a parallel work, Neilson (2006-b) entertains a very similar notion of shame. The questions and the methodology of the two works are different. Section 6 comments in more detail.

knowledge of the menu.⁴ DM has well-defined preferences over sets of alternatives (menus). Our interpretation of shame as the motivating emotion allows considerations of fairness to impact preferences only through their effect on second-stage choices, where the presence of a fairer option reduces the attractiveness of an allocation. The underlying normative notion of fairness is central to our model, because assumptions on the norm of fairness are indirect assumptions on DM's preferences. Assuming a *particular* norm of fairness is difficult, descriptively as well as normatively. Instead, we pose what we consider minimal normative constraints on fairness.

Our representation results establish a correspondence between DM's norm of fairness and her choice behavior. On the one hand, this illustrates how those minimal constraints on fairness impact choice. On the other hand, the particular norm of fairness used by DM can be elicited from her choice behavior.

1.2. Illustration of Results

Denote a typical menu as $A = \{(a_1, a_2), (b_1, b_2), \dots\}$, where the first and second components in each alternative are, respectively, the private payoff for DM and for the recipient. We pose axioms on DM's preferences over menus that allow us to establish a sequence of representation theorems. To illustrate our results, consider a special case of those representations:

$$U(A) = \max_{(a_1, a_2) \in A} [u(a_1) + \beta \varphi(a_1, a_2)] - \beta \max_{(b_1, b_2) \in A} [\varphi(b_1, b_2)],$$

where u and φ are increasing in all arguments. u is a utility function over private payoffs and $\varphi(a_1, a_2)$ is interpreted as the fairness of the allocation (a_1, a_2) .

Alternatively, if we denote by a^* and b^* the two maximizers above, it can be written as:

$$U(A) = \underbrace{u(a_1^*)}_{\text{value of private payoff}} - \underbrace{\beta(\varphi(b_1^*, b_2^*) - \varphi(a_1^*, a_2^*))}_{\text{shame}}.$$

This representation captures the tension between the impulse to maximize private payoff and the desire to minimize shame from not choosing the fairest alternative within a set. It evaluates a menu by the highest utility an allocation on the menu gets, where this utility depends on the menu itself. The utility function that is used to evaluate allocations is additive and has two distinct components. The first component, $u(a_1)$, gives the value of a

⁴If the exit option is chosen in the aforementioned experiment by Dana et al, as in our setup, the recipient does not observe that there was a dictator, who could have chosen another allocation. In their experiment, the recipient is further unaware that another person was involved at all. It would be interesting to see how informing the recipient that some other person had received \$9 would change the experimental findings. This would correspond to our setup.

degenerate menu (a singleton set) that contains the allocation under consideration. When evaluating degenerate menus, which leave DM with a trivial choice under observation, we assume her to be *selfish*: she prefers one allocation to another if and only if the former gives her a greater private payoff, independent of the recipient's payoff. The second component is "shame." It represents the cost DM incurs when selecting (a_1, a_2) in the face of the fairest available alternative, (b_1^*, b_2^*) .

As shame is evoked whenever this fairest available alternative is not chosen, we can relate choice to a second binary relation "fairer than," which represents DM's private norm of fairness. We assume that DM's private norm of fairness induces a *Fairness Ranking* of all alternatives, which is represented by $\varphi(a_1, a_2)$. We further assume that DM's norm of fairness satisfies *Solvability*, implying that the fairness ranking is never satiated in one player's payoff, and the *Pareto* criterion in payoffs, implying that φ is increasing in all arguments.

In the special case considered here, the shame from choosing (a_1, a_2) in stage two is $\beta(\varphi(b_1^*, b_2^*) - \varphi(a_1, a_2))$. Hence, even alternatives that are not chosen may matter for the value of a set, and larger sets are not necessarily better. To see this, consider the representation above with $u(a_1) = a_1$, $\beta = \frac{1}{2}$ and $\varphi(a_1, a_2) = a_1 a_2$. Compare the sets $\{(10, 1), (4, 3)\}$, $\{(10, 1)\}$ and $\{(4, 3)\}$. Evaluating these sets we find $U\{(10, 1), (4, 3)\} = 9$, $U\{(10, 1)\} = 10$ and $U\{(4, 3)\} = 4$. To permit such a ranking, we assume a version of *Left Betweenness*, which allows smaller sets to be preferred over larger sets. Left Betweenness weakens the Set Betweenness assumption first introduced by Gul and Pesendorfer (2001), henceforth GP. Theorem 1 establishes that our weakest representation, which captures the intuition discussed thus far, is equivalent to the collection of all the above assumptions.

Selfishness leaves no room for altruism. Suppose, however, that only the second stage of the procedure is observed (for example, because DM, as in the classic dictator game, never gets to choose between menus). In this case, our representations might conform with DM behaving as if she had direct interest in the recipient's welfare and had to trade off this altruistic motive with concerns about her private payoff. We argue that it is hard to reconcile such an interpretation with observing any choice reversal in stage two. Thus, when observing stage two in isolation, shame can mimic altruism only if the induced choice ranking is set independent. Theorem 2 establishes that, given the assumptions made so far, an additional separability assumption on preferences over sets, *Consistency*, is equivalent to the existence of such a ranking. In the special case of our representation considered above, the induced choice behavior satisfies *Consistency*. To see this, regroup the terms as follows:

$$U(A) = \underbrace{\max_{(a_1, a_2) \in A} [u(a_1) + \beta \varphi(a_1, a_2)]}_{\text{second stage choice criterion}} - \underbrace{\beta \max_{(b_1, b_2) \in A} [\varphi(b_1, b_2)]}_{\text{effect of fairest alternative}}.$$

We further specify the norm of fairness by assuming that the private payoffs to the two players have *Independent Fairness Contributions*: The fairness contribution of raising one player's payoff can not depend on the level of the other player's payoff. The idea is that interpersonal utility comparisons are infeasible. With this additional assumption, Theorem 3 establishes that there are two utility functions, v_1 and v_2 , evaluated in the payoff to DM and the recipient respectively, such that the value of their product represents the fairness ranking, $\varphi(a_1, a_2) = v_1(a_1)v_2(a_2)$. Thus, the fairest alternative within a set of alternatives can be characterized as the Nash Bargaining Solution (NBS) of an associated game. Because the utility functions used to generate this game are private, so is the norm.⁵ We argue that when based on true selfish utilities, the NBS is a convincing fairness criterion in our context. Those utilities, however, may not be publicly known, especially in anonymous choice situations, and therefore, DM may not be able to base her evaluation on true selfish utilities. Nevertheless, one can assess the descriptive appeal of the representation by asking whether the utilities comprising the norm at least resemble selfish utilities.

Example: Let $u(a_1) = a_1$, $\varphi(a_1, a_2) = v_1(a_1)v_2(a_2) = a_1a_2$ and $\beta = \frac{1}{2}$. This implies that selfish utility u is risk neutral and unbounded, and that the utilities v , which are used to generate the fairness ranking, coincide with u . Shame is half the difference between the Nash-product of the fairest and the chosen alternatives. Reconsider the experiment by Dana et al (2006) mentioned above, with the added constraint that only integer values are possible allocations. The set $A = \{(10, 0) (9, 1) (8, 2), \dots, (0, 10)\}$ then corresponds to the dictator game. It induces the imaginary bargaining game with possible utility-allocations $\{(10, 0) (9, 1) (8, 2), \dots, (0, 10), (0, 0)\}$, where the imaginary disagreement point is $\lim_{(x,y) \rightarrow 0} (v_1^{-1}(x), v_2^{-1}(y)) = (0, 0)$. According to the NBS, $(5, 5)$ would be the outcome of the bargaining game. Its fairness is $5 \cdot 5 = 25$. To trade off shame with selfishness, DM chooses the alternative that maximizes the sum of private utility and fairness, $a_1 + a_1a_2$, which is $(6, 4)$. Its fairness is $6 \cdot 4 = 24$ and the shame incurred by choosing it is $\frac{1}{2}$. Hence $U(A) = 5.5$. From the singleton set $B = \{(9, 0)\}$, which corresponds to the exit option in the experiment, the choice is trivial and $U(B) = 9$. This example illustrates both the tension DM is exposed to when choosing from a large set and the reason why she might prefer a smaller menu.

Finally, Theorem 4 extends the former representations by allowing DM to be responsible for the welfare of many other recipients. This extension is then applied to model a social decision maker who is able to alter the transparency of her policies' consequences. Policies

⁵Therefore, the fairness ranking could also be represented by a different functional, based on different utilities.

create social value, but also have a redistributive component. DM faces a trade-off when choosing the transparency of her policies: More transparency makes it easier for the public to perceive fair choices as such, while less transparency makes it harder for society to detect selfish choices. Shame, therefore, might lead her to implement policies with relatively opaque consequences.

The organization of the paper is as follows: Section 2 presents the basic model and a representation that captures the concepts of fairness and shame. Section 3 isolates a choice criterion from the choice situation. Section 4 further specifies the fairness ranking. Section 5 extends the representation to finitely many other players and suggests an application to a social decision maker. Section 6 points out connections to existing literature and section 7 concludes.

2. The Model

Let K be the set of all finite subsets of \mathbb{R}_+^2 .⁶ Any element $A \in K$ is a finite set of alternatives. A typical alternative $\mathbf{a} = (a_1, a_2)$ is interpreted as a payoff pair, where a_1 is the private payoff for DM and a_2 is the private payoff allocated to the (potentially anonymous) other player, the recipient. Endow K with the topology generated by the Hausdorff metric, which is defined for any pair of non-empty sets, $A, B \in K$, by:

$$d_h(A, B) := \max \left[\max_{\mathbf{a} \in A} \min_{\mathbf{b} \in B} d(\mathbf{a}, \mathbf{b}), \max_{\mathbf{b} \in B} \min_{\mathbf{a} \in A} d(\mathbf{a}, \mathbf{b}) \right],$$

where $d : \mathbb{R}_+^2 \rightarrow \mathbb{R}_+$ is the standard Euclidian distance.

Let \succ be a continuous preference relation (weak order) over K . We write $A \succ B$ if DM strictly prefers A to B . The associate weak preference, \succeq and the indifference relation, \sim are defined in the usual way.

The choice of a menu $A \in K$ is not observed by the recipient, while the choice *from* any menu is. We call the impact this observation has on choice "shame." Of course various other regarding preferences that are not impacted by observation could be present as well. We do not account for those, as our aim is not to describe a range of possible attitudes toward others, but to derive a tractable representation according to which DM distinguishes the two stages in an intuitive way.

The first axiom specifies DM's preferences over singleton sets.

P_1 (Selfishness) $\{\mathbf{a}\} \succ \{\mathbf{b}\}$ if and only if $a_1 > b_1$.

⁶With \mathbb{R}_+ we denote the positive reals including 0. \mathbb{R}_{++} denotes the positive reals without 0.

A singleton set $\{\mathbf{a}\}$ is a degenerate menu that contains only one feasible allocation, (a_1, a_2) . It leaves DM with a trivial choice to be made when being observed in the second stage. Therefore, the ranking over singleton sets can be thought of as the ranking over allocations that are imposed on DM. We contend that there is no room for shame in this situation; choosing between two singleton sets reveals DM's "true" preferences over allocation outcomes. The axiom states that DM is not concerned about the payoff to the second player when evaluating such sets; she compares any pair of alternatives based solely on the first component, her private payoff. If, for example, DM had an altruistic concern for fairness in the dictator game previously described, she would strictly prefer the menu $\{(9, 1)\}$ to $\{(9, 0)\}$. P_1 rules out such altruistic concerns. Negative emotions regarding the other player, such as spite or envy, are ruled out as well.

The next axiom captures that shame is a mental cost, which is invoked by unchosen alternatives.

P_2 (Strong Left Betweenness) *If $A \succeq B$, then $A \succeq A \cup B$. Further, if $A \succ B$ and $\exists C$ such that $A \cup C \succ A \cup B \cup C$, then $A \succ A \cup B$.*

We assume that adding unchosen alternatives to a set can only increase shame. Therefore, no alternative is more appealing when chosen from $A \cup B$, than when chosen from one of the smaller sets, A or B . Hence, $A \succeq B$ implies $A \succeq A \cup B$.⁷ Furthermore, if additional alternatives add to the shame incurred by the original choice from a menu $A \cup C$, then they must also add to the shame incurred by any choice from the smaller menu A . Thus, if there is C such that $A \cup C \succ A \cup B \cup C$ and if $A \succ B$, then $A \succ A \cup B$.

Shame, which is the only motive DM knows beyond selfishness, must refer to some personal norm that determines what the appropriate choice should have been. In our interpretation, this norm is to choose one of the fairest available allocations. Interpreting "fairness" as a property of an allocation, which is independent of the menu it is on, we consider a binary relation \succ_f over \mathbb{R}_+^2 as a second primitive.

Definition: If $\mathbf{b} \succ_f \mathbf{a}$, we say that DM considers \mathbf{b} to be *fairer than* \mathbf{a} .

Some of the axioms below are imposed on \succ_f rather than on \succ and are labeled by F instead of P . The underlying notion of fairness is at the heart of those assumptions.⁸ To

⁷This is the "Left Betweenness" axiom. It appears in Dekel, Lipman and Rustichini (2005) and is a weakening of "Set Betweenness" as first posed in GP.

⁸In everyday language, "fair" is sometimes used to capture various different notions. According to the

make them descriptively intuitive, we emphasize their normative appeal, implying that DM will want her norm of fairness to satisfy them. Making these assumptions directly on \succ_f is natural. The relation \succ_f is not directly observable, but the next axiom relates it to observable choice behavior. One contribution of our work is that the implications of F -axioms on \succ are most easily understood from the representation.

P_3 (**Shame**) *If $\exists A \in K$ with $\mathbf{a} \in A$, such that $A \succ A \cup \{\mathbf{b}\}$, then $\mathbf{b} \succ_f \mathbf{a}$.*⁹

$A \succ A \cup \{\mathbf{b}\}$ implies that \mathbf{b} adds to the shame incurred by the original choice in A . The interpretation is that DM is concerned about not choosing a fairest available alternative. Thus, \mathbf{b} must be fairer than any alternative in A , in particular $\mathbf{b} \succ_f \mathbf{a}$.

Definition: We say that DM is *susceptible to shame* if there exists A and B with $A \succ A \cup B$.

Note that for a DM who is purely selfish,¹⁰ \succ_f is empty.

F_1 (**Fairness Ranking**) *\succ_f is an anti-symmetric and negatively transitive binary relation.*

Our discussion rests on the assumption that DM can rank alternatives according to their fairness. In \mathbb{R}_+^2 and with increasing utility from self-payoffs, this assumption is not unreasonably restrictive.¹¹

Combined with P_3 , F_1 implies that only one alternative in each menu, the fairest, is responsible for shame.

F_2 (**Pareto**) *If DM is susceptible to shame, then $\mathbf{a} \geq \mathbf{b}$ and $\mathbf{a} \neq \mathbf{b}$ imply $\mathbf{a} \succ_f \mathbf{b}$.*

According to this axiom, absolute, as opposed to relative, well-being matters; the Pareto

Merriam-Webster Collegiate Dictionary (Tenth Edition, 2001) "*Fair implies an elimination of one's own feelings, prejudices, and desires so as to achieve a proper balance of conflicting interests.*" This is the definition of "fair" we base our arguments on.

⁹The notion of "fairer than" is analogous to the definition of "more tempting than" in Gul and Pesendorfer (2005).

¹⁰This corresponds to the "standard" economic agent, whose preferences satisfy the following variant of P_1 :

$$A \succ B \Leftrightarrow \max_{\mathbf{a} \in A} a_1 > \max_{\mathbf{b} \in B} b_1.$$

Thus, for a purely selfish DM, $A \succ B$ implies $A \sim A \cup B$.

¹¹If, instead, there were a globally most preferred self-payoff, this assumption would rule out very reasonable preference rankings.

criterion excludes notions such as "strict inequality aversion." The resulting concept of fairness must have some concern for efficiency. In the case where there truly is no potential for redistribution, we believe that people find the Pareto criterion a reasonable requirement for one allocation to be fairer than another.¹²

F_3 (**Solvability**) *If $(a_1, 0) \not\sim_f (b_1, b_2)$ then $\exists x$ such that $(a_1, x) \sim_f (b_1, b_2)$. Analogously, if $(0, a_2) \not\sim_f (b_1, b_2)$ then $\exists y$ such that $(y, a_2) \sim_f (b_1, b_2)$.*

Ignoring the qualifier, the axiom states that in order to make two allocations deemed equally fair, any variation in the level of one person's payoff can always be compensated by appropriate variation in the level of the other person's payoff. This requires \succ_f never to be satiated in any person's payoff. Relying on F_2 , the qualifiers take into account that monetary payoffs are bounded below by 0. For example, F_3 implies that there is a sum x , such that $(x, 1) \sim_f (10, 10)$. This assumption captures the insight that any fairness ranking with a concern for efficiency must go beyond the Pareto principle and trade off, in some manner, payoffs across individuals.

As \succ is continuous, \succ_f is continuous in all alternatives for which P_3 relates \succ to \succ_f . $F_1 - F_3$ imply that this is the case on $\mathbb{R}_+ \times \mathbb{R}_{++}$.¹³ Assuming that \succ_f is continuous even in alternatives for which P_3 does not relate \succ to \succ_f has obviously no implication for choice. For ease of exposition, we assume in all what follows that \succ_f is continuous on all of \mathbb{R}_+^2 .

Theorem 1 *If DM is susceptible to shame, then \succ and \succ_f satisfy $P_1 - P_3$ and $F_1 - F_3$ respectively, if and only if there exist continuous and strictly increasing functions $u : \mathbb{R}_+ \rightarrow \mathbb{R}$, $\varphi : \mathbb{R}_+^2 \rightarrow \mathbb{R}$ and a continuous function $g : \mathbb{R}_+^2 \times \varphi(\mathbb{R}_+^2) \rightarrow \mathbb{R}$, weakly increasing in its second argument and satisfying: $g(\mathbf{a}, x) \geq 0$ whenever $\varphi(\mathbf{a}) \leq x$, such that the function $U : K \rightarrow \mathbb{R}$ defined as $U(A) = \max_{\mathbf{a} \in A} \left[u(a_1) - g\left(\mathbf{a}, \max_{\mathbf{b} \in A} \varphi(\mathbf{b})\right) \right]$ represents \succ and φ represents \succ_f . If DM is not susceptible to shame, $g \equiv 0$.*

¹²In many contexts, people would disagree with the statement that the allocation (1million, 6) is fairer than (5, 5). On the basis of the definition in footnote 10, however, we claim that the opposition to (1million, 6) as a fair allocation can only be based on the implicit premise that there must be some mechanism to divide the gains more evenly (Such a mechanism would imply the availability of a third option, which would render both of the above allocations unfair.) In an explicit choice situation this premise cannot be sustained. The Pareto property has indeed been advocated in the philosophical literature on fairness. Rawls (1971), for example, proposes the idea of "original position," a mental exercise whereby a group of rational people must establish a principle of fairness (e.g. when distributing income) without knowing beforehand where on the resulting pecking order they will end up themselves. Requiring that the allocation satisfy *Pareto* makes much sense in such an environment.

¹³ \succ_f is relevant for choice in alternative \mathbf{b} , if and only if there is \mathbf{c} with $\mathbf{c} \prec_f \mathbf{b}$ and $c_1 > b_1$, which requires $c_2 < b_2$. Thus $b_2 > 0$ is necessary for the construction of \mathbf{c} .

All detailed proofs are in the appendix. We now highlight the important steps. As both \succ and \succ_f are continuous binary relations, they can be represented by continuous functions $U : K \rightarrow \mathbb{R}$ and $\varphi : \mathbb{R}_+^2 \rightarrow \mathbb{R}$ respectively. φ is an increasing function as implied by *Pareto* (F_2). The combination of *Strong Left Betweenness* (P_2), *Shame* (P_3) and *Fairness Ranking* (F_1) implies GP's *Set Betweenness* (SB) property: $A \succeq B$ implies $A \succeq A \cup B \succeq B$. GP demonstrate that imposing SB on preferences over sets makes every set indifferent to a certain subset of it, which includes at most two elements (Lemma 2 in their paper). Hence we confine our attention to a subset of our domain that includes all sets with cardinality no greater than 2. *Selfishness* (P_1) and P_3 imply that a set $\{\mathbf{a}, \mathbf{b}\}$ is strictly inferior to $\{\mathbf{a}\}$ if and only if $a_1 > b_1$ and $\mathbf{b} \succ_f \mathbf{a}$. We can then strengthen GP's Lemma 2 and state that any set is indifferent to some two-element set that includes one of the fairest allocations in the original (larger) set. Using *Solvability* (F_3) we show the continuity of the second component, the function g , in the representation.

The representation in Theorem 1 highlights the basic trade-off between private payoff and shame as the only concepts DM may care about. There are at most two essential alternatives within a set, to be interpreted as the "chosen" and the "fairest" alternative, \mathbf{a} and \mathbf{b} respectively. For the latter, its fairness, $\varphi(\mathbf{b})$, is a sufficient statistic for its impact on the set's value. DM suffers from shame, measured by $g(\mathbf{a}, \varphi(\mathbf{b}))$, whenever $\varphi(\mathbf{a}) < \varphi(\mathbf{b})$, where $\varphi(\mathbf{a})$ is the fairness of the chosen alternative. The representation captures the idea of shame being an emotional cost that emerges whenever the fairest available allocation is not chosen. Its magnitude may depend on the fairness of the chosen allocation.

The main contribution of Theorem 1 is the provision of a way to elicit DM's fairness ranking, \succ_f , from choice behavior: all functions in the representation are continuous and hence, for $\mathbf{b} \in \mathbb{R}_+ \times \mathbb{R}_{++}$ and $\mathbf{b} \succ_f \mathbf{a}$, there is \mathbf{c} , such that $U(\{a, c\}) > U(\{a, b, c\})$. Since it is continuous, φ is then uniquely determined on its entire domain, \mathbb{R}_+^2 .

Note that the properties of the function g and the max operator inside imply that the second term is always a cost (non-positive). The other max operator implies that DM's payoff will never lie below b_1 , which is her payoff as suggested by the fairest allocation. Thus, any deviations by DM from choosing the fairest allocation will be in her own favor. These observations justify labeling said cost as "shame."

From the representation, it is easy to see that the induced choice correspondence,

$$C(A) := \left\{ \arg \max_{\mathbf{a} \in A} \left[u(a_1) - g \left(\mathbf{a}, \max_{\mathbf{b} \in A} \varphi(\mathbf{b}) \right) \right] \right\}$$

may be context dependent in the sense that a higher degree of shame may affect choice. In

other words, if we define a binary relation "better choice than," \succ_c , by $\mathbf{a} \succ_c \mathbf{b}$ if $\exists B$ with $\mathbf{b} \in B$, such that $B \cup \{\mathbf{a}\} \succ B$, then this binary relation need not be acyclic. This feature may be plausible when shame is taken into account. In the next section we spell out the implications of enforcing a context-independent criterion for choice.

3. A Second-Stage Choice Ranking

In many situations, only second-stage choice may be observable. For example, the standard dictator game corresponds only to second-stage choice in our setup. Typical behavior in various versions of this game, where subjects tend to give part of the endowment to the recipient, is often interpreted as motivated by an altruistic motive. We interpret altruism to imply that the recipient's welfare is a good, just as selfishness implies that DM's private payoff is a good.¹⁴ If DM had those two motives, she would have to make a trade-off between them. As in the case of two generic goods, very basic assumptions would lead to a context-independent choice ranking of alternatives. As we point out at the end of section 2, we can define a binary relation "better choice than," \succ_c , by $\mathbf{a} \succ_c \mathbf{b}$ if $\exists B$ with $\mathbf{b} \in B$, such that $B \cup \{\mathbf{a}\} \succ B$. This binary relation need not be acyclic: Different choice problems, A and B , may lead to different second-stage rankings of \mathbf{a} and \mathbf{b} , for $\mathbf{a}, \mathbf{b} \in A \cap B$. If no cycles occur, second-stage behavior might look as if it were generated by, for instance, a trade-off of selfishness and altruism, even though observation of stage-one choice would rule this out. If, on the other hand, cycles are observed in stage-two choice, simple altruistic motives cannot be solely responsible for behavior that is not purely selfish. In this section we identify a condition on preferences that makes DM's second-stage choice independent of the choice set. This implies finding a function $\psi : \mathbb{R}_+^2 \rightarrow \mathbb{R}$ that assigns a value to each $\mathbf{a} \in A$, such that \mathbf{a} is a choice from A only if $\psi(\mathbf{a}) \geq \psi(\mathbf{b})$ for all $\mathbf{b} \in A$.

Definition: $X := \{(\mathbf{a}, \mathbf{b}) : \{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}\}$ is the set of all pairs of alternatives generating strict *Set Betweenness*.

For any set of two allocations $\{\mathbf{a}, \mathbf{b}\}$, we interpret the preference ordering $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$ as an indication of a discrepancy between what DM chooses (\mathbf{a}) and the alternative she deems to be the fairest (\mathbf{b}), which causes her choice to bear shame. This shame, however, is not enough to make her choose \mathbf{b} .

Combined with F_1 , Shame (P_3) implies that choice between sets depends on the fairness

¹⁴This interpretation is based on the following definition of altruism (Merriam-Webster Collegiate Dictionary [Tenth Edition, 2001]): "*Unselfish regard for or devotion to the welfare of others.*" We understand this definition as ruling out any considerations that condition on available but unchosen alternatives.

of the fairest alternative in the set. The next axiom relates choice to the fairness of the chosen alternative as well: The fairer DM's choice, the less shame she feels.

P_4 (Fairer is Better) *If for $\{\mathbf{a}\} \sim \{\mathbf{a}'\}$ we have $\{(\mathbf{a}, \mathbf{b}), (\mathbf{a}', \mathbf{b})\} \subseteq X$ and $\mathbf{a} \succ_f \mathbf{a}'$, then $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}', \mathbf{b}\}$.*

Axiom P_4 implies that only the fairness of the chosen alternative matters for its impact on shame.

Given $P_1 - P_4$ and $F_1 - F_3$, an additional separability assumption is equivalent to separable shame, and thus to a set-independent choice ranking.

P_5 (Consistency) *If*

$$\{(\mathbf{a}, \mathbf{b}), (\mathbf{a}, \mathbf{d}), (\mathbf{a}', \mathbf{b}'), (\mathbf{a}', \mathbf{d}'), (\mathbf{c}, \mathbf{b}), (\mathbf{c}', \mathbf{b}'), (\mathbf{c}, \mathbf{d}), (\mathbf{c}', \mathbf{d}')\} \subseteq X,$$

then $\{\mathbf{a}, \mathbf{b}\} \sim \{\mathbf{a}', \mathbf{b}'\}$ and $\{\mathbf{a}, \mathbf{d}\} \sim \{\mathbf{a}', \mathbf{d}'\}$ imply $\{\mathbf{c}, \mathbf{b}\} \succ \{\mathbf{c}', \mathbf{b}'\} \Leftrightarrow \{\mathbf{c}, \mathbf{d}\} \succ \{\mathbf{c}', \mathbf{d}'\}$.

We make no claim about the normative or descriptive appeal of this assumption. Instead, we view it as an empirical criterion: If the condition is not met, observation of stage-two choice should suffice to distinguish altruism from shame as the motive behind DM's other-regarding behavior. The axiom requires independence between the impact of the chosen and the fairest alternative on the set ranking:

$$\{(\mathbf{a}, \mathbf{b}), (\mathbf{a}, \mathbf{d}), (\mathbf{a}', \mathbf{b}'), (\mathbf{a}', \mathbf{d}'), (\mathbf{c}, \mathbf{b}), (\mathbf{c}', \mathbf{b}'), (\mathbf{c}, \mathbf{d}), (\mathbf{c}', \mathbf{d}')\} \subseteq X$$

implies that from each of the sets $\{\mathbf{a}, \mathbf{b}\}$, $\{\mathbf{a}, \mathbf{d}\}$, $\{\mathbf{a}', \mathbf{b}'\}$, $\{\mathbf{a}', \mathbf{d}'\}$, $\{\mathbf{c}, \mathbf{b}\}$, $\{\mathbf{c}', \mathbf{b}'\}$, $\{\mathbf{c}, \mathbf{d}\}$ and $\{\mathbf{c}', \mathbf{d}'\}$, the alternative listed first is chosen in the second stage despite the availability of a fairer alternative, which is listed second. Assume, without loss of generality that $\{\mathbf{a}\} \succ \{\mathbf{a}'\}$. Suppose there are two pairs of fairer and less attractive alternatives, \mathbf{b}, \mathbf{b}' and \mathbf{d}, \mathbf{d}' , such that for each of them pairing their members with \mathbf{a} and \mathbf{a}' , respectively, gives rise to indifference. In the context of Theorem 1, this implies that both pairs induce the same shame differential, which exactly cancels the selfish preference of $\{\mathbf{a}\}$ over $\{\mathbf{a}'\}$: $\{\mathbf{a}, \mathbf{b}\} \sim \{\mathbf{a}', \mathbf{b}'\}$ and $\{\mathbf{a}, \mathbf{d}\} \sim \{\mathbf{a}', \mathbf{d}'\}$. Then, the axiom states that pairing the members of \mathbf{b}, \mathbf{b}' or \mathbf{d}, \mathbf{d}' with any other chosen alternatives \mathbf{c} and \mathbf{c}' , respectively, must also lead to the same differential in shame. In particular, $\{\mathbf{c}, \mathbf{b}\} \succ \{\mathbf{c}', \mathbf{b}'\}$ implies $\{\mathbf{c}, \mathbf{d}\} \succ \{\mathbf{c}', \mathbf{d}'\}$. Again, the validity of this technical assumption in a given context is an empirical question.

Theorem 2 *If DM is susceptible to shame, then \succ and \succ_f satisfy P_1-P_5 and F_1-F_3 respectively, if and only if there exist continuous and strictly increasing functions $u : \mathbb{R}_+ \rightarrow \mathbb{R}$ and $\varphi : \mathbb{R}_+^2 \rightarrow \mathbb{R}$, such that the function $U : K \rightarrow \mathbb{R}$ defined as $U(A) = \max_{\mathbf{a} \in A} [u(a_1) + \varphi(a_1, a_2)] - \max_{\mathbf{b} \in A} [\varphi(b_1, b_2)]$ represents \succ and φ represents \succ_f .
If DM is not susceptible to shame, h is a constant.*

The proof constructs a path in the (a_1, a_2) -plane such that the fairness $\varphi(\mathbf{a})$ increases along this path. Then, on two neighboring indifference curves in the $(\mathbf{a}, \varphi(\mathbf{b}))$ -space, $\varphi(\mathbf{b})$ increases, as \mathbf{a} varies along the path. Relying on P_5 , these indifference curves allow us to rescale $\varphi(\mathbf{b})$ to make the representation of \succ quasi-linear.¹⁵ Separability is then immediate. Since the proof of Theorem 2 is a special case of the proof of Theorem 4, we only go through the more general case in detail in the appendix.

The representation isolates a choice criterion that is independent of the choice problem: DM's behavior is governed by maximizing

$$u(a_1) + \varphi(a_1, a_2).$$

The value of the set is reduced by

$$\max_{\mathbf{b} \in A} \varphi(b_1, b_2),$$

a term that depends solely on the fairest alternative in the set. Grouping the terms differently reveals the trade-off between self-payoff, $u(a_1)$, and the shame involved with choosing \mathbf{a} from the set A :

$$\max_{\mathbf{b} \in A} [\varphi(b_1, b_2) - \varphi(a_1, a_2)] \geq 0.$$

Note that now shame takes an additively separable form, depends only on the fairness of both alternatives, and is increasing in the fairness of the fairest and decreasing in that of the chosen alternative. If $P_1 - P_4$ and $F_1 - F_3$ hold, then, according to Theorem 2, P_5 is equivalent to having a set-independent choice ranking.

4. Specifying a Fairness Ranking

In this section we impose one more axiom on \succ_f to further characterize the fairness ranking. It asserts that the fairness contribution of one person's marginal payoff cannot depend on the initial payoff levels.

¹⁵A more elaborate discussion on this technique appears after Theorem 3.

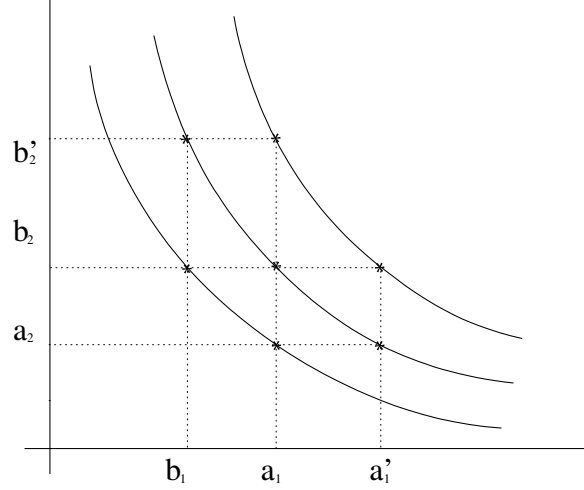


Figure 1: Independent Fairness Contributions.

F_4 (Independent Fairness Contributions) If $(a_1, a_2) \sim_f (b_1, b_2)$ and $(a'_1, a_2) \sim_f (a_1, b_2) \sim_f (b_1, b'_2)$, then $(a'_1, b_2) \sim_f (a_1, b'_2)$.

The axiom is illustrated in figure 1. If $a_1 = a'_1$ or $b_2 = b'_2$, this axiom is implied by F_1 , F_2 and the continuity of \succ_f . For $a_1 \neq a'_1$ and $b_2 \neq b'_2$, the statement is more subtle. Consider first a stronger assumption:

F'_4 (Strong Independent Fairness Contributions) $(a_1, a_2) \sim_f (b_1, b_2)$ and $(a'_1, a_2) \sim_f (b_1, b'_2)$ imply $(a'_1, b_2) \sim_f (a_1, b'_2)$.

The fairness contribution of one person's marginal payoff cannot depend on the initial payoff level of the other person: It is unclear to DM how much an increase in monetary payoff means to the recipient, because even if the (marginal) utility of the recipient were known to DM, she could not compare it to her own, as interpersonal utility comparisons are infeasible. The qualifier in F'_4 establishes that DM considers the fairness contribution of changing her own payoff from a_1 to a'_1 given the allocation (a_1, a_2) to be the same as that of changing the recipient's payoff from b_2 to b'_2 given (b_1, b_2) . F'_4 then states that starting from the allocation (a_1, b_2) , changing a_1 to a'_1 should again be as favorable in terms of fairness as changing b_2 to b'_2 . This is the essence of *Independent Fairness Contributions*. The stronger qualifier $(b_1, b'_2) \sim_f (a_1, b_2) \sim_f (a'_1, a_2)$ in F_4 weakens the axiom. For example, the fairness ranking $(a_1, a_2) \succ_f (b_1, b_2)$ if and only if $\min(a_1, a_2) > \min(b_1, b_2)$ is permissible under F_4 ,

but not under F'_4 .¹⁶

Theorem 3 \succ_f satisfies $F_1 - F_4$, if and only if there are continuous, increasing and unbounded functions $v_1, v_2 : \mathbb{R}_+ \rightarrow \mathbb{R}_{++}$, such that $\varphi(\mathbf{a}) = v_1(a_1)v_2(a_2)$ represents \succ_f .

Luce and Tukey (1964) prove the necessity and sufficiency of *Solvability* and the *Corresponding Trade-offs Condition* (the label they use for F_4) to admit an additive representation. To show how a proof works, we repeatedly use axiom F_4 to establish that if $(a_1, a_2) \sim_f (a'_1, a'_2)$ and $(a_1, \tilde{a}_2) \sim_f (a'_1, \tilde{a}'_2)$, then $(\tilde{a}_1, a_2) \sim_f (\tilde{a}'_1, a'_2) \Leftrightarrow (\tilde{a}_1, \tilde{a}_2) \sim_f (\tilde{a}'_1, \tilde{a}'_2)$. With this knowledge, we can create a monotone increasing mapping $a_2 \rightarrow \gamma(a_2)$ that transforms the original indifference map to be quasi-linear with respect to the first coordinate in the $(a_1, \gamma(a_2))$ plane. Keeney and Raiffa (1976) refer to the procedure we employ as the lock-step procedure. Quasi-linearity implies that there is an increasing continuous function $\xi : \mathbb{R}_+ \rightarrow \mathbb{R}$, such that $\varphi(\mathbf{a}) := \xi(a_1) + \gamma(a_2)$ represents \succ_f . Define $v_1(a_1) := \exp(\xi(a_1))$ and $v_2(a_2) := \exp(\gamma(a_2))$. Then $v_1, v_2 : \mathbb{R}_+ \rightarrow \mathbb{R}_{++}$ are increasing and continuous and if we redefine $\varphi(\mathbf{a}) := v_1(a_1)v_2(a_2)$, it represents \succ_f .

This representation suggests an appealing interpretation of the fairness ranking DM is concerned about: She behaves *as if* she had in mind two increasing and unbounded utility functions, one for herself¹⁷ and one for the recipient. By mapping the alternatives within each set into the associated utility space, any choice set induces a finite bargaining game where only the disagreement point is unspecified. DM then identifies the fairest alternative within a set *as if* she also had in mind a disagreement point, that makes this alternative the Nash Bargaining Solution¹⁸ of the game.¹⁹ Moreover, the fairness of all alternatives can be ranked according to the same functional, namely the Nash product.

Remember that F_3 requires trading off marginal payoffs. The tension of having to trade off marginal payoffs without being able to compare their welfare contribution (F_4) is common in a range of social-choice problems.²⁰ Our axioms are weak in the sense that they do not

¹⁶ F_4 is referred to as the *Hexagon condition* or the *Corresponding Trade-offs Condition* (Keeney and Raiffa [1976]), F'_4 as the *Thomsen condition*. With F_2 and F_3 , F'_4 is implied by F_4 . See Karni and Safra (1998) for a proof.

¹⁷This utility function need not agree with her true utility for personal payoffs, u . The interpretation is that DM is concerned about the recipient's perception of her choice. The recipient, however, may not know DM's true utility, especially under anonymity.

¹⁸See Nash (1950).

¹⁹The imaginary disagreement point is determined by $\lim_{(x,y) \rightarrow 0} (v_1^{-1}(x), v_2^{-1}(y))$. It could be some finite and weakly positive pair of utility payoffs. In particular it could be $(0, 0)$, which corresponds to DM imagining that players walk away in the case that no agreement is reached. It could also be negative. This corresponds to DM imagining that players have an extra incentive to find an agreement: there is a cost to disagreement.

²⁰For a review, see Hammond (1990).

constrain DM in this trade-off, as long as she takes into account that the fairness contribution of increasing one person's payoff should not depend on the other's payoff. The power of Theorem 3 is that it bases a representation on these weak assumptions. The downside is that the form of this representation is not unique, as the utilities v_1 and v_2 are not observable independent of the norm of fairness. For example, there is another pair of increasing utility functions such that DM is concerned about their sum, that is, she acknowledges efficiency as the only fairness criterion.

To underline the appeal of the Nash product as a descriptive representation of fairness,²¹ we now point out how DM might reason within the constraints of the axioms:

We justified the Pareto criterion, F_2 , as a plausible axiom for the fairness ranking. As argued above, concern for fairness requires the acknowledgment of some form of interpersonal comparability of preferences' intensity. If utilities were known cardinally, symmetry in terms of utility payoffs is the other criterion we would expect the ranking to satisfy.²² In our context, this implies independence of the role people play, dictator or recipient. However, utilities are inherently ordinal, rendering such a comparison infeasible. At best we can, if we assume people to have cardinal utilities that reflect their attitudes toward risk, determine marginal utilities up to scaling. Mariotti (1997), for example, considers a context in which "*interpersonal comparisons of utility are meaningful; that is, there exists an (unknown) rescaling of each person's utility which makes utilities interpersonally comparable.*" At the same time, however, "*interpersonal comparisons of utility are not feasible.*" Assume there is a correct interpersonal utility scaling, but DM cannot determine it. Can she guarantee that for this unknown scaling both symmetry and Pareto are satisfied? They would have to be satisfied for all potential scalings. Mariotti establishes that the NBS is the only criterion with this property.

Even more appealing is an interpretation of the NBS as the fairest allocation that is related to Gauthier's (1986) principle of "moral by agreement": Trying to assess what is fair, but finding herself unable to compare utilities across individuals, DM might refer to the prediction of a symmetric mechanism for generating allocations. In particular, DM might ask what would be the allocation if both she and the recipient were to bargain over the division of the surplus. To answer this question, she does not need to assume the intensities of the two preferences. This is a procedural interpretation that is not built on the axioms: DM is not ashamed of payoffs, but of using her stronger position in distributing the gains. It is,

²¹Even though u and v_1 do not have to agree, our interpretation might be more convincing when they resemble each other empirically. In particular it is more appealing if DM's actual utility from self-payoff u is unbounded.

²²This reasoning leads Rawls (1971) to suggest *Pareto* and *Symmetry* as the two criteria a decision maker under a veil of ignorance should respect.

then, the intuitive and possibly descriptive appeal of the NBS in many bargaining situations that makes it normatively appealing to DM in our context.²³ Theorem 3 establishes the behavioral equivalence of this interpretation and our axioms.

The Pareto and the Solvability axioms, F_2 and F_3 respectively, rule out fairness rankings with $(x, 0) \sim_f (0, y)$ for all x, y . In particular the Nash product with linear utility functions v_1, v_2 is ruled out as a criterion for fairness. Such orderings could easily be accommodated by posing *Pareto* and *Solvability* only on \mathbb{R}_{++}^2 . As a consequence, φ would be strictly increasing only on \mathbb{R}_{++}^2 and v_1, v_2 would only have to be weakly positive, $v_1, v_2 : \mathbb{R}_+ \rightarrow \mathbb{R}_+$.²⁴ These weaker axioms would still rule out the Maximin as a criterion for fairness.

Remark: Any concern DM has about fairness originates from being observed. Consequently, DM should expect a potentially anonymous observer to share her notion of what is fair: Her private norm of fairness, which we observe indirectly, should reflect her concern about not violating a social norm. If the observed choice situation is anonymous, DM does not know the recipient's identity and is aware that the recipient does not know hers. Therefore, the ranking cannot depend on either identity. Combining this with the idea that fairness of an allocation should not depend on the role a person plays, whether dictator or recipient, one might want to pose symmetry of the fairness ranking in terms of direct payoffs.

F_5 (**Symmetry**) $(a_1, a_2) \sim_f (a_2, a_1)$.

Adding this assumption constrains $v_1(a) = v_2(a)$ in the representation of Theorem 3. The numerical example given in the introduction features the combination of Theorem 2 and Theorem 3, where all functions involved are the identity. For brevity, we will not repeat it here.

5. Multiple Recipients and an Application

In order to expand the range of possible applications of our representation, we first extend our results to finitely many agents.

²³The descriptive value of the NBS has been tested empirically. For a discussion see Davis and Holt (1993) pages 247-55. Further, multiple seemingly natural implementations of it have been proposed (Nash [1953], Osborne and Rubinstein [1994]).

²⁴As can be seen in the proof of Theorem 2, this would imply the possibility of $(-\infty, -\infty)$ as an imaginary disagreement point, which corresponds to DM imagining that players have to find an agreement (infinite cost of disagreement).

5.1. Multiple Recipients

The underlying idea is that DM (without loss of generality individual 1) is concerned about $N - 1 \geq 2$ other individuals, whose payoffs depend on her choice. In analogy to section 2, let K be the set of all finite subsets of \mathbb{R}_+^N . Any element $A \in K$ is a finite set of alternatives. A typical alternative $\mathbf{a} = (a_1, a_2, \dots, a_N)$ is interpreted as a payoff vector, where a_n is the payoff allocated to individual n . We write, for example, $(a_m, a_n, \mathbf{a}_{-m,n})$ as the alternative with payoff a_m to individual m , payoff a_n to individual n and $\mathbf{a}_{-m,n} \in \mathbb{R}_+^{N-2}$ lists all other individuals' payoffs in order. We endow K with the topology generated by the Hausdorff metric.

Let \succ be a continuous preference relation over K . Most of the axioms we pose on \succ in section 2 can be readily applied to \succ on this new domain. We define \succ_f in analogy to the previous definition. Instead of F_3 we write

F_3^N (**Weak Solvability**) *If $(a_n, \mathbf{0}) \not\succeq_f \mathbf{b}$ then for all $m \neq n$, there exists a_m such that $(a_m, a_n, \mathbf{0}) \sim_f \mathbf{b}$.*

The axiom states that it is always possible to equate the fairness of an allocation with payoff to only one individual to that of an initially fairer allocation by giving appropriate payoffs to any second individual. This property requires the fairness ranking never to be satiated in any individual payoff.

Definition: The pair of possible payoffs to individuals m and n is *Preferentially Independent with respect to its Complement* (P.I.C.), if the fairness ranking in the (a_m, a_n) -space is independent of $\mathbf{a}_{-m,n}$.

F_4^N (**Pairwise Preferential Independence**) *For all $m, n \in \{1, \dots, N\}$, the pair of possible payoffs to individuals m and n is P.I.C.*

Similarly to F_4 , this axiom must hold if the contribution of one person's marginal private payoff to the fairness of an allocation cannot depend on another person's private payoff level.

Theorem 4 *Assume $N \geq 3$ and that DM is susceptible to shame.*

(i) \succ and \succ_f satisfy $P_1 - P_5$ and F_1, F_2 and F_3^N respectively, if and only if there exist continuous and strictly increasing functions $u : \mathbb{R}_+ \rightarrow \mathbb{R}$ and $\varphi : \mathbb{R}_+^N \rightarrow \mathbb{R}$ such that the function $U : K \rightarrow \mathbb{R}$ defined as $U(A) = \max_{\mathbf{a} \in A} [u(a_1) + \varphi(a_1, a_2, \dots, a_n)] - \max_{\mathbf{b} \in A} [\varphi(b_1, b_2, \dots, b_n)]$ represents \succ and φ represents \succ_f .

(ii) \succ_f also satisfies F_4^N if and only if there exist continuous and strictly increasing functions $v_1, \dots, v_N : \mathbb{R}_+ \rightarrow \mathbb{R}_{++}$, where v_1, \dots, v_N are unbounded such that $\varphi(\mathbf{a}) = \prod_{i=1}^N v_i(a_i)$.
If DM is not susceptible to shame, φ is constant.

Theorem 4 is analogous to Theorem 2. For the proof, note that the analogue of Theorem 1 can be established by substituting \mathbf{a}_{-1} for a_2 in the theorem and in the proof, where now $\varphi : \mathbb{R}_+^N \rightarrow \mathbb{R}$. To establish the analogue of Theorem 3, namely that there are N increasing unbounded functions v_1, \dots, v_N , such that the fairness ranking \succ_f can be represented by $\varphi(\mathbf{a}) = \prod_{i=1}^N v_i(a_i)$ if and only if it satisfies F_1, F_2, F_3^N and F_4^N , we first state a stronger version of F_3^N :

$F_3^{N'}$ (**Solvability**) If $(a_n, \mathbf{a}_{-n}) \not\sim_f \mathbf{b}$ then for all $m \neq n$, there exists a_m such that $(a_m, a_n, \mathbf{a}_{-m,n}) \sim_f \mathbf{b}$.

We observe that *Continuity*, F_1 , F_2 and F_3^N imply *Solvability*. To see this, assume $(a_n, \mathbf{a}_{-n}) \not\sim_f \mathbf{b}$. By F_2 , $(a_n, \mathbf{0}) \not\sim_f (a_n, \mathbf{a}_{-n})$ and hence (using F_1) $(a_n, \mathbf{0}) \not\sim_f \mathbf{b}$. By F_3^N , there exists \tilde{a}_m such that $(\tilde{a}_m, a_n, \mathbf{0}) \sim_f \mathbf{b}$. By F_2 again, $(\tilde{a}_m, a_n, \mathbf{z}) \succeq_f \mathbf{b}$ for all $\mathbf{z} \in \mathbb{R}_+^{N-2}$. Therefore, by *Continuity*, there must be $a_m \in \mathbb{R}_+$ for which $(a_m, a_n, \mathbf{a}_{-m,n}) \sim_f \mathbf{b}$. We can then apply:

Theorem (Luce and Tukey [1964]) *Pairwise Preferential Independence and Solvability imply the existence of an additive representation of \succ_f .*

The proof of this theorem can be found in Krantz et al (1971). We illustrate the idea for the case $N = 3$ by showing that F_4^N implies F_4 for (without loss of generality) the pair of individuals 1 and 2, independent of the payoff to individual 3:

For any (a_1^0, a_2^0, a_3^0) and any a_1^1 , define a_2^1 and a_3^1 such that

$$(a_1^1, a_2^0, a_3^0) \sim_f (a_1^0, a_2^1, a_3^0) \sim_f (a_1^0, a_2^0, a_3^1).$$

Applying F_4^N twice implies that

$$(a_1^1, a_2^1, a_3^0) \sim_f (a_1^1, a_2^0, a_3^1) \sim_f (a_1^0, a_2^1, a_3^1).$$

For any a_1^2 , define a_2^2 and a_3^2 such that

$$(a_1^2, a_2^0, a_3^0) \sim_f (a_1^0, a_2^2, a_3^0) \sim_f (a_1^0, a_2^0, a_3^2) \sim_f (a_1^1, a_2^1, a_3^0).$$

We have to show that $(a_1^2, a_2^1, a_3) \sim_f (a_1^1, a_2^2, a_3)$ for any value of a_3 : $(a_1^2, a_2^0, a_3^0) \sim_f (a_1^1, a_2^0, a_3^1)$, so by F_4^N also $(a_1^2, a_2^1, a_3^0) \sim_f (a_1^1, a_2^1, a_3^1)$. Similarly $(a_1^0, a_2^2, a_3^0) \sim_f (a_1^0, a_2^1, a_3^1)$, so by F_4^N also $(a_1^1, a_2^2, a_3^0) \sim_f (a_1^1, a_2^1, a_3^1)$. Using transitivity, $(a_1^2, a_2^1, a_3^0) \sim_f (a_1^1, a_2^2, a_3^0)$ and by F_4^N this is independent of a_3^0 . Hence $(a_1^2, a_2^1, a_3) \sim_f (a_1^1, a_2^2, a_3)$ for any value of a_3 .

The existence of utility functions according to which \succ_f is represented by the Nash product follows, as before, where additivity is implied by Luce and Tukey's theorem. We gave the intuition for the remainder of the proof of Theorem 4 after stating Theorem 2.

5.2. An Application to Obfuscation by a Social Decision Maker

It is often argued that individuals who make social choices are faced with very rigid constraints. Shame at acting against the interests of others could be one such constraint, moderating individuals' decisions as compared to their selfish interest. We build on this interpretation to explain why a social decision maker may implement policies with relatively opaque consequences. To first illustrate by simple example why such lack of transparency (or obfuscation) might be valuable to DM at all, consider an indivisible good that can be assigned to one individual. All individuals have the same probability of needing it the most. Under obfuscation, this uncertainty never gets resolved, hence all allocations are equally fair and DM can take the good for herself without shame. If, on the other hand, the uncertainty does get resolved before DM chooses an allocation, she can only claim the good without shame in the event that she values it the most.

The literature that studies obfuscation in policy making usually considers redistributive policies. As an example, Tullock (1983) uses the decision of where to locate a new road: Depending on the road's location, some citizens will gain, others might lose. These consequences will not be entirely transparent at the time of decision making.

While building a road in a certain location clearly has a redistributive component, we assert that it may also generate value for the society as a whole. In this section, we therefore consider more general policies, which carry both an uncertain social value and an uncertain distribution of gains among citizens. All citizens (including DM) have identical information with respect to both types of uncertainty at every stage of the process.²⁵ Evaluating policies requires some degree of public deliberation, for example, the consultation of experts. Before

²⁵This assumption stands in contrast to the usual asymmetric-information assumption (either among citizens or between citizens and DM), that is used to explain the choice among different methods of redistribution. See, for example, Coate and Morris (1995) and a survey in Wittman (1989).

this deliberation takes place, DM can limit the degree to which deliberation will resolve uncertainty. She does so by choosing the transparency level of the policies that will be considered.²⁶ We assume that even for the lowest feasible transparency level, deliberation will reveal DM's selfish payoff from each relevant policy. This assumption is intuitive due to DM's arguably exposed role. It is also appropriate when addressing DM, who is constrained by shame. Finally, it is crucial for the results established below, as the assumption introduces an asymmetry between DM and all other citizens despite the information structure assumed above: The probability that DM's preferences will become public is larger than that of any other citizen's. Therefore, when deciding on the transparency of policies, DM has to trade off the benefit of obfuscation, which makes selfish choices seem more fair, and the value of transparency, which reveals efficient choices as such.

The time sequence is as follows: Firstly, DM chooses the transparency level for the policies under consideration. Secondly, public deliberation symmetrically reduces the uncertainty about the consequences of those policies. The higher the transparency level was set, the less uncertainty remains. Lastly, DM chooses one policy. It is important to note that, in slight contrast to our model, stage-one choice does not alter, in terms of expected payoffs, the set of policies that are relevant for stage-two choice. Instead, it alters the expected differential in fairness between the policies in which DM will have a selfish interest, and those that will be perceived as fairest.

Formally, consider a very large population of N individuals indexed by i . Let $\Omega = \mathbb{R}$ be identified both as the set of possible policy choices $a \in \Omega$ and as the type space. When referring to a particular individual i , we denote her type as $x_i \in \Omega$. Individual i 's selfish preferences are commonly known to be represented by $u : \mathbb{R}^+ \rightarrow \mathbb{R}^+$, which is continuous and strictly decreasing with the (standard Euclidian) distance, $d(a, x_i)$, between the implemented policy a and her type x_i . Types are identically and independently distributed according to a Normal distribution with an unknown mean, $\theta \in \mathbb{R}$, and known variance $\sigma^2 > 0$. The (conjugate) common prior distribution of θ is Normal, $\theta \sim N(\bar{\theta}_0, \bar{\nu}_0^2)$, where $\bar{\nu}_0^2 := \nu^2$ and, without loss of generality, $\bar{\theta}_0 = 0$. Thus, σ captures the uncertainty about the redistributive consequences of policies, while ν relates to the uncertainty about the value generated for society as a whole. Let Ω^n be the set of all possible type profiles of length $n \leq N$, with a typical element $x^n = (x_1, \dots, x_n)$. Upon observing the realization $x^n \in \Omega^n$, each individual, including DM, updates her beliefs according to Bayes' rule. The resulting common posterior distribution of θ is Normal as well, $\theta \sim N(\bar{\theta}_n, \bar{\nu}_n^2)$, with $\bar{\theta}_n = \frac{\sigma^2(x_1 + \dots + x_n)}{\nu^2 + n\sigma^2}$ and $\bar{\nu}_n^2 = \frac{\nu^2\sigma^2}{\nu^2 + n\sigma^2}$.

²⁶For example, DM can set the agenda of issues she wants to address: Instead of debating the location of the road, she could also choose to deliberate introducing a tax. The individual consequences of the tax are presumably more transparent than those of the location of the road.

Note that $x_1 \in x^n$. Since we identify DM as individual 1, her type is always revealed for $n \geq 1$.

After the entire population observes x^n , DM makes a social choice $a \in \Omega$. We assume that DM's preferences satisfy $P_1 - P_5$ as well as F_1, F_2, F_3^N and F_4^N . As is implied by Theorem 4, DM would like to choose a to maximize

$$u(d(a, x_1)) + \beta_N \prod_{i=1}^N u(d(a, x_i)) - \max_{b \in \Omega} \left[\beta_N \prod_{i=1}^N u(d(b, x_i)) \right],$$

where, for simplicity, we assume h to be linear ($h(x) = \beta_N x$). After observing x^n , however, x_{n+1}, \dots, x_N remain unknown, so $\prod_{i=1}^N u(d(a, x_i))$ cannot be evaluated. To accommodate the uncertainty about the distribution of types in the population, we assume instead that the *expected* fairness of an allocation conditional on x^n ,

$$E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right]$$

is what determines shame. Thus, the relevant characteristics of a policy are both its proximity $d(a, x_1)$ to DM's type x_1 and its expected fairness. According to our representation, the action that the public perceives as the fairest is

$$a^* := \arg \max_{a \in \Omega} E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right].$$

DM's choice is then governed by maximizing the term

$$u(d(a, x_1)) + \beta_N \left(E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right] - E \left[\prod_{i=1}^N u(d(a^*, x_i)) \mid x^n \right] \right).$$

Note that for fixed n and large population size, $N \rightarrow \infty$, $a^* \rightarrow \bar{\theta}_n$.

Before x^n is observed by both DM and the public, DM can pick $n \in [1, \bar{n}]$, $\bar{n} \ll N$. The number n is interpreted as the transparency level of policies in Ω : The more transparent the policies are, the larger the number of individuals whose type becomes revealed by the public deliberation. Different transparency levels n introduce different distributions over expected fairness, while leaving $d(a, x_i)$ unchanged. Thus, the choice of n is equivalent to stage-one choice of different distributions over menus with the same cardinality that differ in the expected fairness of their elements. In contrast to the policy choice, we assume that the choice of the transparency level, n , is free of shame.²⁷ This assumption could rest on the fact

²⁷Due to this assumption, the model nicely fits our general framework.

that the transparency level is chosen before any uncertainty is resolved and, hence, cannot bias the ex post expected fairness of any policy. Alternatively, the public might simply be unaware of transparency as a choice dimension.

We are interested in finding the optimal transparency level, that is, the optimal first-stage choice according to DM's selfish preferences. DM faces the following trade-off: On the one hand, she benefits from high transparency, which reduces the uncertainty about the fairness of policies and allows the public to interpret fair choices as such. On the other hand, she benefits from low transparency, as it gives her selfish payoff more weight in public observation, limiting the public's ability to detect selfish behavior.

To determine DM's optimal transparency choice, n^* , define the ratio of the standard deviations σ and ν as $s := \frac{\sigma}{\nu}$ and let $s = s^*$ solve $2 + 3s^2 - 3s^4 - 6s^6 - 2s^8 = 0$, $s^* \approx 0.84$. Ignoring the integer constraint on n , we state:

Proposition $n^*(s)$ exists and is unique. For $s < s^*$, $m = n^*(s)$ is the solution to $2 + (2m + 1)s^2 - 3s^4 - 2(2m + 1)s^6 - m(m + 1)s^8 = 0$, which is decreasing in s . For $s \geq s^*$, $n^*(s) = 1$.

Note that the optimal transparency level does not depend on DM's susceptibility to shame, β_N .²⁸ This means that while for the case of a standard economic agent with $\beta_N = 0$, the choice of n is arbitrary, even a small positive β_N implies the same unique amount of obfuscation as an arbitrarily large β_N does. Note as well that absolute uncertainty is irrelevant for the optimal transparency choice, only relative uncertainty $s = \frac{\sigma}{\nu}$ matters. This makes the prediction of the proposition very robust.

The proof of the proposition is in the appendix. It establishes that DM's utility is decreasing in the absolute value of the random variable $X_n := \bar{\theta}_n - x_1$, and that X_n is normally distributed. Then $\frac{p(|X_n|=z)}{p(|X_m|=z)}$ can be shown to satisfy the Monotone Likelihood Ratio Property (MLRP). Since DM's utility is decreasing in z , she strictly prefers m over n if and only if $\frac{p(|X_n|=z)}{p(|X_m|=z)}$ is increasing in z . Assuming $n > m$, we find, with some straightforward algebra, that this is the case if and only if $2 + (m + n)s^2 - 3s^4 - 2(m + n)s^6 - mns^8 < 0$. Thus DM has well-defined preferences over levels of transparency, n , and these preferences depend only on s . We then establish that $n^*(s)$ is unique and for $s < s^*$ is also decreasing. As a result, if DM prefers $n = 1$ to $n = 2$, then $n^*(s) = 1$ is her globally preferred transparency level. If she prefers n to both $n - 1$ and $n + 1$, then $n^*(s) = n$ is her globally preferred level.

²⁸The proposition is only concerned with the transparency choice. The allocation DM chooses in the second stage obviously does depend on β_N , as it determines the extent to which DM yields to shame at the cost of her selfish interest.

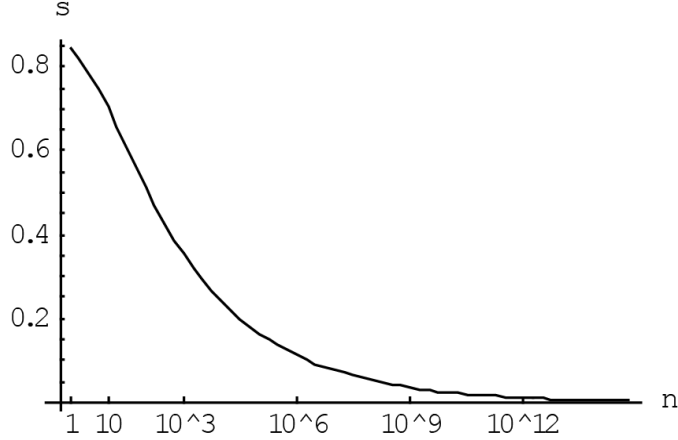


Figure 2: n as a function of s .

Before interpreting this proposition, consider again the trade-off DM faces: Since her utility is decreasing in the random variable $|X_n| := |\bar{\theta}_n - x_1|$, she would like the mean of the public posterior on types to be as close as possible to her type. Since DM's type is always observed (and thus always affects the posterior), it has a higher expected weight in the posterior than any other citizen's, which is only observed with probability $\frac{n-1}{N}$. The lower n is, the greater the advantage DM has over other citizens. Lowering n , however, increases the weight that the common prior gets in the public posterior. Since types are correlated, this is not in DM's interest. Now consider the proposition in the context of this trade-off: DM prefers to consider more opaque policies for the future if the standard deviation of the distribution of benefits across the population, measured by σ , becomes larger compared to the uncertainty about which policy is socially optimal, which is measured by ν . Intuitively, in this case she is concerned about the situation where her selfish preferences conflict with considerations of fairness. Therefore, she would like her own selfish preferences to impact the public posterior as much as possible. Since her preferences always become public, she would like other citizens' preferences to remain unobserved. Conversely, she prefers more transparent policies if the uncertainty is more about the socially optimal policy and less about the distribution of gains. In this case, she is mostly concerned about the situation where her selfish preferences are in line with considerations of fairness but an uninformed public would perceive her as selfish if she chose accordingly. Figure 2 shows $n^*(s)$.

6. Related Literature

Other-regarding preferences have been considered extensively in economic literature. In particular, inequality aversion as studied by Fehr and Schmidt (1999) is based on an objective function with a similar structure to the representation of second-stage choice in Theorem 3.²⁹ Both works attach a cost to any deviation from choosing the fairest alternative. In Fehr and Schmidt's work, the fairest allocation need not be feasible and is independent of the choice situation. In our work, the fairest allocation is always a feasible choice and it is identified through the axioms. This dependence of the fairest allocation on the choice situation allows us to distinguish observed from unobserved choice.

The idea that there may be a discrepancy between DM's preference to behave "pro-socially" and her desire to be viewed as behaving pro-socially is not new to economic literature. For a model thereof, see Benabou and Tirole (2006).

Neilson's (2006-b) work is motivated by the same experimental evidence as ours. He also considers menus of allocations as objects of choice. Neilson does not axiomatize a representation result, but points out how choices among menus should relate to choices from menus, if shame were the relevant motive. He relates the two aspects of shame that also underlie the *Set Betweenness* property in our work; DM might prefer a smaller menu over a larger menu either because avoiding shame compels her to be generous when choosing from the larger menu, or because being selfish when choosing from the larger menu bears the cost of shame.

The structure of our representation resembles the representation of preferences with self-control under temptation, as axiomatized in GP. GP study preferences over sets of lotteries and show that their axioms lead to a representation of the following form:

$$U^{GP}(A) = \max_{a \in A} \{u^{GP}(a) + v^{GP}(a)\} - \max_{b \in A} \{v^{GP}(b)\}$$

with u^{GP} and v^{GP} both linear in the probabilities and where A is now a set of lotteries. In their context, u^{GP} represents the "commitment"- and v^{GP} the "temptation"-ranking. While the two works yield representations with a similar structure, their domains - and therefore the axioms - are different. In particular, the objects in GP's work are sets of lotteries. They impose the independence axiom and indifference to the timing of the resolution of uncertainty. This allows them to identify the representation above that consists of two functions that are linear in the probabilities. Each of these functions is an expected utility functional. The objects in our work, in contrast, are sets of allocations and there is no uncertainty. The

²⁹Neilson (2006-a) axiomatizes a reference-dependent preference, that can be interpreted in terms of Fehr and Schmidt's objective function.

natural way to introduce uncertainty to our model is to treat our representation as the utility function, which should be used to calculate the expected utility of lotteries over sets. Therefore, DM would typically not be indifferent to the timing of the resolution of uncertainty.³⁰ However, one of GP's axioms is the *Set Betweenness* axiom, $A \succ B \Rightarrow A \succ A \cup B \succ B$. We show that our axioms *Strong Left Betweenness* (P_2), *Shame* (P_3) and *Fairness Ranking* (F_1) imply *Set Betweenness*. Hence, GP's Lemma 2 can be employed, allowing us to confine attention to sets with only two elements.

Our model is positive in nature, but it is interesting to contrast moral or normative elements in its interpretation with those in the context of the temptation literature: In a work related to GP, Dekel, Lipman and Rustichini (2005) write: "...by 'temptation' we mean that the agent has some view of what is normatively correct, what she should do, but has other, conflicting desires which must be reconciled with the normative view in some fashion." According to this interpretation, the commitment ranking is given a normative value. In our work, shame is based on deviating from some fairness norm that tells DM what she should do. This norm conflicts with DM's selfish commitment ranking.

Empirically, the assumption that only two elements of a choice set matter for the magnitude of shame (the fairest available alternative and the chosen alternative) is clearly simplifying: Oberholzer-Gee and Eichenberger (2004) observe that dictators choose to make much smaller transfers when their choice set includes an unattractive lottery. In other words, the availability of an unattractive allocation seems to lessen the incentive to share.

Lastly, it is necessary to qualify our leading example: The experimental evidence on the effect of (anonymous) observation on the level of giving in dictator games is by no means conclusive. Behavior tends to depend crucially on surroundings, like the social proximity of the group of subjects and the phrasing of the instructions, as, for example, Bolton, Katok and Zwick (1994); Burnham (2003); and Haley and Fessler (2005) record. On the one hand, there is more evidence in favor of our interpretation: In a follow-up to the experiment cited in the introduction, Dana et al (2006) verify that dictators do not exit the game if second-stage choice is also unobserved. Similarly, Pillutla and Murningham (1995) find evidence that people's giving behavior under anonymity depends on the information given to the observing recipient. In experiments related to our leading example, Lazear, Malmendier and Weber (2005) as well as Broberg, Ellingsen and Johannesson (2008) even predict and find that the most generous dictators are keenest to avoid an environment where they could share with an observing recipient.³¹ Broberg et al further elicit the price subjects are willing to pay in order to exit the dictator game, finding that the mean exit reservation price equals

³⁰In section 5.2 we account for uncertainty, which can be translated into uncertainty over sets.

³¹This nicely underlines our interpretation of "shame" as a motive.

82% of the dictator game endowment. On the other hand, this is in contrast to evidence collected by Koch and Normann (2005), who claim that altruistic behavior persists at an almost unchanged level when observability is credibly reduced. Similarly, Johannesson and Persson (2000) find that incomplete anonymity - not observability - is what keeps people from being selfish. Ultimately, experiments aimed at eliciting a norm share the same problem: Since people use different (and potentially contradictory) norms in different contexts, it is unclear whether the laboratory environment triggers a different set of norms than would other situations: Frohlich, Oppenheimer and Moore (2000) point out that money might become a measure of success rather than a direct asset in the competition-like laboratory environment, such that the norm might be "do well" rather than "do not be selfish."³² More theoretically, Miller (1999) suggests that the phrasing of instructions might determine which norm is invoked. For example, the reason that Koch and Normann do not find an effect of observability might be that their thorough explanation of anonymity induces a change in the regime of norms, in effect telling people "be rational," which might be interpreted as "be selfish." Then being observed might have no effect on people who, under different circumstances, might have been ashamed to be selfish.

7. Conclusion

We study a decision maker who cares about others' well-being only when being observed. We term the motive that distinguishes choice behavior when observed from choice when unobserved "shame." Theorem 1 features a representation that captures the tension between the interest to maximize private payoff and the shame that results from not choosing the fairest alternative within a set. Theorem 2 identifies a set-independent choice criterion with the help of a separability axiom. If there is a set-independent choice criterion, the representation should be more tractable for applications. More importantly, the separability assumption provides a criterion on preferences over sets to decide whether or not the period-two choice of alternatives might look as if it was generated by an altruistic concern. In Theorem 3 we further specify the norm of fairness. We show that the fairest alternative in a set can be interpreted as the Nash Bargaining Solution of an associated game. Because the utility functions used to generate this game are private, so is the norm. The most appealing interpretation relies on the descriptive power of the NBS in many bargaining situations, giving it normative appeal as the solution to a symmetric mechanism. Lastly, Theorem 4 extends Theorem 2 to situations where DM faces multiple other individuals whose welfare

³²Surely the opposite is also conceivable: Subjects might be particularly keen to be selfless when the experimenter observes their behavior. This example is just ment to draw attention to the difficulties faced by experimenters in the context of norms.

depends on her choice. We apply our model to a social decision maker, whose selfish utility is correlated with fairness. She faces a trade-off when choosing the transparency of her policies: Being more transparent makes it easier for the public to perceive fair choices as such, while less transparency makes it harder for society to detect selfish choices. In our setup, the optimal transparency level is unique and is independent of DM's susceptibility to shame.

Let us conclude with another experiment to suggest how to incorporate uncertainty into our model. Dana, Weber and Kuang (2005) make a dictator face a choice between \$5 and \$6 for herself. An anonymous recipient will receive either \$5 or \$1. Which recipient payoff is paired with which dictator payoff is determined by a coin flip. The dictator can reveal (without being observed) the outcome of the coin flip prior to her decision. The authors find that many dictators choose not to reveal the outcome. This action seems weakly dominated, because whether or not the dictator is willing to give up \$1 in order to give the recipient the extra \$4, knowing whether such a trade-off is necessary should not hurt DM. We propose to interpret the revealed and the unrevealed conditions as two different choice situations. If all functions in the combination of Theorem 2 and Theorem 3 are identities, and if DM subscribes to the vNM axioms, the utilities to be compared are

$$U(\{(6, 3), (5, 3)\}) = 6$$

versus

$$\frac{1}{2}U(\{(6, 5), (5, 1)\}) + \frac{1}{2}U(\{(6, 1), (5, 5)\}) = \frac{1}{2}6 + \frac{1}{2}5 = 5.5$$

This could explain the observed behavior. However, since in the experiment the recipient knows the full instructions and does not observe DM's decision to reveal, observability would require a more involved interpretation. To ease the application of our model, it would be interesting to see how the subject's behavior changes if the recipient is only told the information DM has *after* her decision to reveal or not.

8. Appendix

8.1. Proof of Theorem 1

Let $U : K \rightarrow \mathbb{R}$ be a continuous function that represents \succ . Define $u(a_1) \equiv U(\{(a_1, 0)\})$. By P_1 , $u(a_1) = U(\{(a_1, a_2)\})$ independent of a_2 , with $u(a_1)$ continuous and strictly increasing.

Let $\varphi : \mathbb{R}_+^2 \rightarrow \mathbb{R}$ be a continuous function that represents \succ_f . By F_2 , φ is also strictly increasing.

Claim 1.1 (Right Betweenness): $A \succeq B \Rightarrow A \cup B \succeq B$.

Proof: There are two cases to consider:

Case 1) $\forall a \in A, \exists b \in B$ such that $b \succ_f a$. Let $A = \{\mathbf{a}^1, \mathbf{a}^2, \dots, \mathbf{a}^N\}$ and $C_0 = B$. Define $C_n = C_{n-1} \cup \{\mathbf{a}^n\}$ for $n = 1, 2, \dots, N$. According to F_1 , there exists $b \in B$ such that $a^n \not\succeq_f b$. By P_3 , $C_{n-1} \not\succeq C_n$. By negative transitivity of \succ , $C_0 \not\succeq C_N$ or $A \cup B \succeq B$.

Case 2) $\exists a \in A$ such that $a \succ_f b, \forall b \in B$. Let $B = \{\mathbf{b}^1, \mathbf{b}^2, \dots, \mathbf{b}^M\}$. Define $C_0 = A$ and $C_m = C_{m-1} \cup \{\mathbf{b}^m\}$ for $m = 1, 2, \dots, M$. By P_3 , $\forall C$ such that $a \in C, C \not\succeq C \cup \{\mathbf{b}^m\}$. Hence, $C_{m-1} \not\succeq C_m$. By negative transitivity of \succ , $C_0 \not\succeq C_M$ or $A \cup B \succeq A \succeq B$, hence $A \cup B \succeq B$. \parallel

Combining Claim 1.1 with P_2 guarantees Set Betweenness (SB): $A \succeq B \Rightarrow A \succeq A \cup B \succeq B$. Having established Set Betweenness, we can apply GP Lemma 2, which states that any set is indifferent to a specific two-element subset of it.

Lemma 1.1 (GP Lemma 2): *If \succ satisfies SB, then for any finite set A , there exist $\mathbf{a}, \mathbf{b} \in A$ such that $A \sim \{\mathbf{a}, \mathbf{b}\}$, (\mathbf{a}, \mathbf{b}) solves $\max_{\mathbf{a}' \in A} \min_{\mathbf{b}' \in A} U(\{\mathbf{a}', \mathbf{b}'\})$ and (\mathbf{b}, \mathbf{a}) solves $\min_{\mathbf{b}' \in A} \max_{\mathbf{a}' \in A} U(\{\mathbf{a}', \mathbf{b}'\})$.*

Define $f : \mathbb{R}_+^2 \times \mathbb{R}_+^2 \rightarrow \mathbb{R}$ such that $f(\mathbf{a}, \mathbf{b}) = u(a_1) - \tilde{U}(\mathbf{a}, \mathbf{b})$, where $\tilde{U} : \mathbb{R}_+^2 \times \mathbb{R}_+^2 \rightarrow \mathbb{R}$ is a function satisfying:

$$U(\{\mathbf{a}, \mathbf{b}\}) = \max_{\mathbf{a}' \in \{\mathbf{a}, \mathbf{b}\}} \min_{\mathbf{b}' \in \{\mathbf{a}, \mathbf{b}\}} \tilde{U}(\mathbf{a}', \mathbf{b}') = \min_{\mathbf{b}' \in \{\mathbf{a}, \mathbf{b}\}} \max_{\mathbf{a}' \in \{\mathbf{a}, \mathbf{b}\}} \tilde{U}(\mathbf{a}', \mathbf{b}').^{33}$$

By definition we have $f(\mathbf{a}, \mathbf{a}) = 0$ for every $\mathbf{a} \in \mathbb{R}_+^2$. Note as well that

$$\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \Rightarrow f(\mathbf{a}, \mathbf{b}) > 0,$$

as otherwise we would have:

$$U(\{\mathbf{a}, \mathbf{b}\}) = \max \left\{ \begin{array}{l} u(a_1) - \max_{f(\mathbf{a}, \mathbf{b})} \{f(\mathbf{a}, \mathbf{a})=0\} \\ u(b_1) - \max_{f(\mathbf{b}, \mathbf{b})} \{f(\mathbf{b}, \mathbf{a})\} \end{array} \right\} \geq u(a_1) - \max \left\{ \begin{array}{l} f(\mathbf{a}, \mathbf{a}) = 0 \\ f(\mathbf{a}, \mathbf{b}) \end{array} \right\} = U(\{\mathbf{a}\}).$$

For a decision maker who is not susceptible to shame, $U(\{\mathbf{a}, \mathbf{b}\}) = \max\{u(a_1), u(b_1)\}$. Hence setting $f(\mathbf{a}, \mathbf{b}) \equiv 0$ is consistent with her preferences. The following claims help us to further identify f for a decision maker who is susceptible to shame.

³³Note that $\max_{\mathbf{a} \in A} \min_{\mathbf{b} \in A} U(\{\mathbf{a}, \mathbf{b}\}) = \max_{\mathbf{a} \in A} \min_{\mathbf{b} \in A} \left[\max_{\mathbf{a}' \in \{\mathbf{a}, \mathbf{b}\}} \min_{\mathbf{b}' \in \{\mathbf{a}, \mathbf{b}\}} \tilde{U}(\mathbf{a}', \mathbf{b}') \right] = \max_{\mathbf{a} \in A} \min_{\mathbf{b} \in A} \tilde{U}(\mathbf{a}, \mathbf{b})$.

Claim 1.2: (i) $[\varphi(\mathbf{a}) < \varphi(\mathbf{b}) \text{ and } a_1 > b_1] \Leftrightarrow \{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\}$

(ii) $[\varphi(\mathbf{a}) < \varphi(\mathbf{b}) \text{ and } a_1 \leq b_1] \Rightarrow \{\mathbf{a}\} \sim \{\mathbf{a}, \mathbf{b}\}$

(iii) $[\varphi(\mathbf{a}) = \varphi(\mathbf{b}) \text{ and } a_1 > b_1] \Rightarrow \{\mathbf{a}\} \sim \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$.

Proof: (i) If $\varphi(\mathbf{b}) > \varphi(\mathbf{a})$ then there exists A such that $\mathbf{a} \in A$ and $A \succ A \cup \{\mathbf{b}\}$. As $a_1 > b_1 \Leftrightarrow \{\mathbf{a}\} \succ \{\mathbf{b}\}$, by P_2 $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\}$. Conversely if $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\}$, then $\mathbf{b} \succ_f \mathbf{a}$ and hence $\varphi(\mathbf{a}) < \varphi(\mathbf{b})$. Further from SB and P_1 , $a_1 > b_1$.

(ii) If $a_1 \leq b_1$ then by SB $\{\mathbf{b}\} \succeq \{\mathbf{a}, \mathbf{b}\}$. Since $\varphi(\mathbf{b}) > \varphi(\mathbf{a})$, there is no B such that $\mathbf{b} \in B$ and $B \succ B \cup \{\mathbf{a}\}$, hence $\{\mathbf{b}\} \sim \{\mathbf{a}, \mathbf{b}\}$.

(iii) By P_1 $\{\mathbf{a}\} \succ \{\mathbf{b}\}$ and SB $\{\mathbf{a}\} \succeq \{\mathbf{a}, \mathbf{b}\}$. As $\varphi(\mathbf{a}) = \varphi(\mathbf{b})$, using (i) we have $\{\mathbf{a}\} \sim \{\mathbf{a}, \mathbf{b}\}$.||

Let $(\mathbf{a}^*(A), \mathbf{b}^*(A))$ be the solution of

$$\max_{\mathbf{a}' \in A} \min_{\mathbf{b}' \in A} U(\{\mathbf{a}', \mathbf{b}'\})$$

so $(\mathbf{b}^*(A), \mathbf{a}^*(A))$ solves $\min_{\mathbf{b}' \in A} \max_{\mathbf{a}' \in A} U(\{\mathbf{a}', \mathbf{b}'\})$.

Claim 1.3: There exists $\mathbf{b} \in \arg \max_{\mathbf{a}' \in A} \varphi(\mathbf{a}')$ such that $A \sim \{\mathbf{a}', \mathbf{b}\}$ for some $\mathbf{a}' \in A$ and $\mathbf{b}^*(A) = \mathbf{b}$.

Proof: Assume not, then there exist \mathbf{a}, \mathbf{c} such that $A \sim \{\mathbf{a}, \mathbf{c}\}$, $(\mathbf{a}, \mathbf{c}) = (\mathbf{a}^*(A), \mathbf{b}^*(A))$. Therefore,

$$\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{c}\} \sim \{\mathbf{a}, \mathbf{b}, \mathbf{c}\} \sim A^{34} \forall \mathbf{b} \in \arg \max_{\mathbf{a}' \in A} \varphi(\mathbf{a}')$$

and hence $\mathbf{c} \succ_f \mathbf{b}$, which is a contradiction.||

For the remainder of the proof, let $I_f(\varphi) := \{\mathbf{b}' : \varphi(\mathbf{b}') = \varphi\}$. That is, $I_f(\varphi(\mathbf{b}))$ is the \sim_f equivalence class of \mathbf{b} . Define

$$Y(\mathbf{a}, \varphi) = \{\mathbf{b}' \in I_f(\varphi) : \{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}'\} \succ \{\mathbf{b}'\}\}$$

We make the following four observations:

(1) $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$, $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{c}\}$ and $\mathbf{b} \succ_f \mathbf{c}$ imply $\{\mathbf{a}, \mathbf{c}\} \succeq \{\mathbf{a}, \mathbf{b}\}$.

(2) $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$, $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{c}\} \succ \{\mathbf{c}\}$ and $\mathbf{b} \sim_f \mathbf{c}$ imply $\{\mathbf{a}, \mathbf{c}\} \sim \{\mathbf{a}, \mathbf{b}\}$.

³⁴Note that if (\mathbf{a}, \mathbf{c}) ((\mathbf{c}, \mathbf{a})) solves the maximin- (minimax-) problem over A , it clearly solves this problem over the subset $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ for all $\mathbf{b} \in A \setminus \{\mathbf{a}, \mathbf{c}\}$.

(3) $\mathbf{b} \in Y(\mathbf{a}, \varphi)$, $\mathbf{b}' \sim_f \mathbf{b}$ and $\{\mathbf{b}\} \succ \{\mathbf{b}'\}$ imply $\mathbf{b}' \in Y(\mathbf{a}, \varphi)$.

(4) If $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$, $\{\mathbf{b}'\} \succ \{\mathbf{b}\}$ and $\mathbf{b}' \in I_f(\varphi(\mathbf{b}))$, then either $\{\mathbf{a}, \mathbf{b}'\} \sim \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}'\}$ or $\{\mathbf{a}, \mathbf{b}'\} \sim \{\mathbf{b}'\} \succeq \{\mathbf{a}, \mathbf{b}\}$.

To verify these observations, suppose first that (1) did not hold. Then $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{c}\}$ and $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$, hence by SB $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ and therefore $\mathbf{c} \succ_f \mathbf{b}$, which is a contradiction. If (2) did not hold, we would get a contradiction to $\mathbf{b} \sim_f \mathbf{c}$ immediately. Next suppose that (3) did not hold. Then $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\} \succ \{\mathbf{b}'\} \sim \{\mathbf{a}, \mathbf{b}'\}$. Note that by SB $\{\mathbf{b}\} \succeq \{\mathbf{b}, \mathbf{b}'\}$ and, applying SB again, $\{\mathbf{b}\} \succeq \{\mathbf{a}, \mathbf{b}, \mathbf{b}'\}$. But then $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{b}, \mathbf{b}'\}$, contradicting $\mathbf{b}' \sim_f \mathbf{b}$. To verify (4), assume $\{\mathbf{a}, \mathbf{b}'\} \succ \{\mathbf{b}'\}$. Then by Claim 1.2 (i) $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}'\} \succ \{\mathbf{b}'\}$ and then by observation (2) $\{\mathbf{a}, \mathbf{b}'\} \sim \{\mathbf{a}, \mathbf{b}\}$. If on the other hand $\{\mathbf{a}, \mathbf{b}'\} \sim \{\mathbf{b}'\}$, then if $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{b}'\}$, $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\}$ and SB imply $\{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{a}, \mathbf{b}, \mathbf{b}'\}$, a contradiction to $\mathbf{b}' \in I_f(\varphi(\mathbf{b}))$. Note that by Claim 1.3 we cannot have $\{\mathbf{b}'\} \succ \{\mathbf{a}, \mathbf{b}'\}$.||

Next we claim that $\varphi(\mathbf{b})$ is a sufficient statistic for the impact of \mathbf{b} on a two element set.

Claim 1.4: There exists a function \tilde{U} satisfying the condition specified above such that $\varphi(\mathbf{b}) > \varphi(\mathbf{a})$ implies $f(\mathbf{a}, \mathbf{b}) = g(\mathbf{a}, \varphi(\mathbf{b}))$, where $g : \mathbb{R}_+^2 \times \mathbb{R} \rightarrow \mathbb{R}$ is weakly increasing in its second argument.

Proof: Such \tilde{U} exists, if and only if $f(\mathbf{a}, \mathbf{b}) = g(\mathbf{a}, \varphi(\mathbf{b}))$ is consistent with \succ . Therefore it is enough to consider the constraints \succ puts on f . Given \mathbf{a} and \mathbf{b} , look at all \mathbf{c} such that $\varphi(\mathbf{b}) > \varphi(\mathbf{c})$. We should show that $f(\mathbf{a}, \mathbf{b}) \geq f(\mathbf{a}, \mathbf{c})$.

First note that if $\varphi(\mathbf{b}) \geq \varphi(\mathbf{a}) \geq \varphi(\mathbf{c})$, then $f(\mathbf{a}, \mathbf{b}) \geq 0 \geq f(\mathbf{a}, \mathbf{c})$ is consistent with \succ . If $\varphi(\mathbf{a}) \geq \varphi(\mathbf{b}) > \varphi(\mathbf{c})$, then $0 \geq f(\mathbf{a}, \mathbf{b}) \geq f(\mathbf{a}, \mathbf{c})$ is consistent with \succ . If $a_1 = 0$, then $f(\mathbf{a}, \mathbf{b}) \geq f(\mathbf{a}, \mathbf{c}) \geq 0$ is consistent with \succ . Therefore, confine attention to the case where $a_1 > 0$ and $\varphi(\mathbf{b}) > \varphi(\mathbf{c}) > \varphi(\mathbf{a})$.

By Claim 1.2 (i), F_2 and F_3 , there exists $\mathbf{b}' \in I_f(\varphi(\mathbf{b}))$ such that $\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}'\}$. Thus, there are two cases to consider:

- 1) $Y(\mathbf{a}, \varphi(\mathbf{b})) \neq \emptyset$.
- 2) $Y(\mathbf{a}, \varphi(\mathbf{b})) = \emptyset$.

Case 1) Suppose $Y(\mathbf{a}, \varphi(\mathbf{b})) \neq \emptyset$. Define $f(\mathbf{a}, \mathbf{b}) := f(\mathbf{a}, \mathbf{b}')$ for some $\mathbf{b}' \in Y(\mathbf{a}, \varphi(\mathbf{b}))$ (note that by observation (2) $f(\mathbf{a}, \mathbf{b}') = f(\mathbf{a}, \mathbf{b}'') \forall \mathbf{b}', \mathbf{b}'' \in Y(\mathbf{a}, \varphi(\mathbf{b}))$ and using observations (3) and (4), this definition is consistent with \succ .) If $Y(\mathbf{a}, \varphi(\mathbf{c})) \neq \emptyset$ then by observation (1) $\{\mathbf{a}, \mathbf{c}\} \succeq \{\mathbf{a}, \mathbf{b}\}$ and hence $f(\mathbf{a}, \mathbf{b}) \geq f(\mathbf{a}, \mathbf{c})$. If $Y(\mathbf{a}, \varphi(\mathbf{c})) = \emptyset$ then $\forall \mathbf{c}' \in I_f(\mathbf{c})$, $\{\mathbf{a}, \mathbf{c}'\} \sim \{\mathbf{c}'\}$. By F_2 and continuity of \succ_f , there exists $\mathbf{c}' \in I_f(\mathbf{c})$ with $c'_1 < b'_1$ for some $\mathbf{b}' \in Y(\mathbf{a}, \varphi(\mathbf{b}))$. Then by Claim 1.1, P_1 and observation (1)

$\{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{c}'\} \succeq \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\} \succ \{\mathbf{c}\}$, so $\mathbf{c}' \in Y(\mathbf{a}, \varphi(\mathbf{c}))$. Contradiction.

Case 2) Suppose $Y(\mathbf{a}, \varphi(\mathbf{b})) = \emptyset$. Define $f(\mathbf{a}, \mathbf{b}) := u(\mathbf{a}_1) - u(0)$. If $Y(\mathbf{a}, \varphi(\mathbf{c})) \neq \emptyset$, then $f(\mathbf{a}, \mathbf{c}) < u(\mathbf{a}_1) = f(\mathbf{a}, \mathbf{b})$. If $Y(\mathbf{a}, \varphi(\mathbf{c})) = \emptyset$ then $f(\mathbf{a}, \mathbf{c}) = u(\mathbf{a}_1) = f(\mathbf{a}, \mathbf{b})$.||

Let $S := \{(\mathbf{a}, \varphi) : Y(\mathbf{a}, \varphi) \neq \emptyset\}$. Note that S is an open set.

Claim 1.5: There is $g(\mathbf{a}, \varphi)$, which is continuous.

Proof: If $Y(\mathbf{a}, \varphi) \neq \emptyset$, then $g(\mathbf{a}, \varphi) = u(\mathbf{a}_1) - U(\{\mathbf{a}, \mathbf{b}\})$ for some $\mathbf{b} \in Y(\mathbf{a}, \varphi)$, and is clearly continuous. If $Y(\mathbf{a}, \varphi) = \emptyset$, then $\varphi \leq \varphi(\mathbf{a})$ implies $g(\mathbf{a}, \varphi) \leq 0$, while $\varphi > \varphi(\mathbf{a})$ implies $g(\mathbf{a}, \varphi) \geq u(\mathbf{a}_1) - u(0)$. Define a switch point $(\widehat{\mathbf{a}}, \widehat{\varphi})$ to be a boundary point of S such that there exists $\widehat{\mathbf{b}} \in \mathbb{R}_+^2$ with $\varphi(\widehat{\mathbf{b}}) = \widehat{\varphi}$. For $\widehat{\varphi} = \varphi(\widehat{\mathbf{a}})$ define $g(\widehat{\mathbf{a}}, \widehat{\varphi}) := 0$ and for $\widehat{\varphi} > \varphi(\widehat{\mathbf{a}})$ define $g(\widehat{\mathbf{a}}, \widehat{\varphi}) := u(\widehat{\mathbf{a}}_1) - u(0)$.

Consider a sequence $\{(\mathbf{a}^n, \varphi^n)\} \rightarrow (\widehat{\mathbf{a}}, \widehat{\varphi})$ in S . Pick a sequence $\{\mathbf{b}^{n'}\}$ with $\mathbf{b}^{n'} \in Y(\mathbf{a}^n, \varphi^n) \forall n$. Define $\{b_1^n\} = \left\{ \min \left[\frac{1}{n}, b_1^{n'}, \widehat{b}_1 \right] \right\}$. Define b_2^n to be a solution to $\varphi(b_1^n, b_2^n) = \varphi^n$. By F_2 and F_3 , b_2^n is well defined. Note that by observation (3) $\mathbf{b}^n = (b_1^n, b_2^n) \in Y(\mathbf{a}^n, \varphi^n)$. Lastly, let $\widehat{b}_1^n \equiv b_1^n$ and \widehat{b}_2^n be the solution to $\varphi(\widehat{b}_1^n, \widehat{b}_2^n) = \widehat{\varphi}$. We have $U(\{\mathbf{a}^n, \mathbf{b}^n\}) = u(\mathbf{a}_1^n) - g(\mathbf{a}^n, \varphi^n)$. If in the switch point $\widehat{\varphi} = \varphi(\widehat{\mathbf{a}})$, then $U(\{\widehat{\mathbf{a}}, \widehat{\mathbf{b}}^n\}) = u(\widehat{\mathbf{a}}_1)$. By continuity, $U(\{\mathbf{a}^n, \mathbf{b}^n\}) - U(\{\widehat{\mathbf{a}}, \widehat{\mathbf{b}}^n\}) \xrightarrow{n \rightarrow \infty} 0$, hence

$$\lim_{n \rightarrow \infty} g(\mathbf{a}^n, \varphi^n) = \lim_{n \rightarrow \infty} [u(\mathbf{a}_1^n) - u(\widehat{\mathbf{a}}_1)] = u(\widehat{\mathbf{a}}_1) - u(\widehat{\mathbf{a}}_1) = 0 = g(\widehat{\mathbf{a}}, \widehat{\varphi}).$$

If in the switch point $\widehat{\varphi} > \varphi(\widehat{\mathbf{a}})$, then $U(\{\widehat{\mathbf{a}}, \widehat{\mathbf{b}}^n\}) = u(\widehat{b}_1^n) = u(b_1^n)$. By the same continuity argument

$$\lim_{n \rightarrow \infty} g(\mathbf{a}^n, \varphi^n) = \lim_{n \rightarrow \infty} [u(\mathbf{a}_1^n) - u(b_1^n)] = u(\widehat{\mathbf{a}}_1) - u(0) = g(\widehat{\mathbf{a}}, \widehat{\varphi}).$$

For $\varphi < \varphi(\mathbf{a})$ let $g(\mathbf{a}, \varphi) < 0$. This satisfies the constraint on f . So g can be continuous in both arguments and increasing in φ and such that for any sequence $\{(\mathbf{a}^n, \varphi^n)\}$ in S , with $\{(\mathbf{a}^n, \varphi^n)\} \rightarrow (\widehat{\mathbf{a}}, \widehat{\varphi})$, we have $\lim_{n \rightarrow \infty} g(\mathbf{a}^n, \varphi^n) = 0$.||

That the representation satisfies the axioms is easy to verify. This completes the proof of Theorem 1.³⁵ ■

³⁵If F_2 and F_3 were only posed on \mathbb{R}_{++}^2 as suggested in section 3, we would have to choose $\widehat{b}_1 > 0$ and $b_1^n > 0$ to use these axioms. This is possible for any switch point other than $(\widehat{\mathbf{a}}, \widehat{\varphi}) = (\mathbf{0}, \varphi(0))$, for which continuity can be established easily.

8.2. Proof of Theorem 2

Theorem 2 and Theorem 4 (i) are analogous, where Theorem 2 covers the case $N = 2$, while Theorem 4 (i) covers the case $N \geq 3$. We prove Theorem 4 (i) below by first establishing that the analogous version of Theorem 1 holds. From there on the proof of Theorem 2 is identical to the proof of Theorem 4 (i), with a_2 substituted for \mathbf{a}_{-1} .

8.3. Proof of Theorem 3

Luce and Tukey [1964] prove the necessity and sufficiency of Solvability (which is implied by Negative Transitivity, Weak Solvability, Pareto and Continuity (apply corollary 1 in the text to the case $N=2$)) and the Corresponding Trade-offs Condition (the label they use for F_4) to admit an additive representation.³⁶ To see how a proof works, consider the Lock-Step Procedure,³⁷ as illustrated by Figure 3:

By F_2 , \succ_f indifference curves are downward sloping and continuous. Fix (a_1^0, a_2^0) and $a_2^1 > a_2^0$. Recursively construct a flight of stairs between the indifference curves through (a_1^0, a_2^0) and (a_1^0, a_2^1) .

In the direction of increasing a_2 (and hence decreasing a_1) :

a_1^n solves $(a_1^n, a_2^n) \sim_f (a_1^0, a_2^0)$. F_3 guarantees that a solution exists whenever $(0, a_2^n) \preceq_f (a_1^0, a_2^0)$. If $(0, a_2^n) \succ_f (a_1^0, a_2^0)$, the flight of stairs terminates.

a_2^{n+1} solves $(a_1^n, a_2^{n+1}) \sim_f (a_1^0, a_2^1)$. A solution exists by F_3 , as $(a_1^n, 0) \prec_f (a_1^0, a_2^1)$ by F_2 .

In the direction of decreasing a_2 (and increasing a_1):

a_1^{-n} solves $(a_1^{-n}, a_2^{-n+1}) \sim_f (a_1^0, a_2^1)$. A solution exists by F_3 , as $(0, a_2^{-n+1}) \prec_f (a_1^0, a_2^1)$ by F_2 .

a_2^{-n} solves $(a_1^{-n}, a_2^{-n}) \sim_f (a_1^0, a_2^0)$. F_3 guarantees that a solution exists whenever $(a_1^{-n}, 0) \preceq_f (a_1^0, a_2^0)$. If $(a_1^{-n}, 0) \succ_f (a_1^0, a_2^0)$, the flight of stairs terminates.

By construction $(a_1^{n+1}, a_2^{n+2}) \sim_f (a_1^n, a_2^{n+1})$ and then by F_4 , $(a_1^n, a_2^{n+2}) \sim_f (a_1^{n-1}, a_2^{n+1})$. Thus we have constructed a discrete set of points on another indifference curve from the initial two curves. Repeating this procedure we can fill \mathbb{R}_+^2 with countable sets of points on countably many indifference curves.

Now consider a particular indifference curve that lies between two members of this set, as illustrated in Figure 4: Define $(a_1^{\frac{1}{2}}, a_2^{\frac{1}{2}})$ implicitly by $(a_1^{\frac{1}{2}}, a_2^{\frac{1}{2}}) \sim_f (a_1^0, a_2^{\frac{1}{2}})$ and $(a_1^{\frac{1}{2}}, a_2^{\frac{1}{2}}) \sim_f (a_1^0, a_2^0)$. Construct a flight of stairs between the indifference curves through $(a_1^0, a_2^{\frac{1}{2}})$ and through (a_1^0, a_2^0) as described above. Then we have in direction of decreasing

³⁶Their theorem is stated in section 5.1 of the text.

³⁷See Keeney and Raiffa (1976).

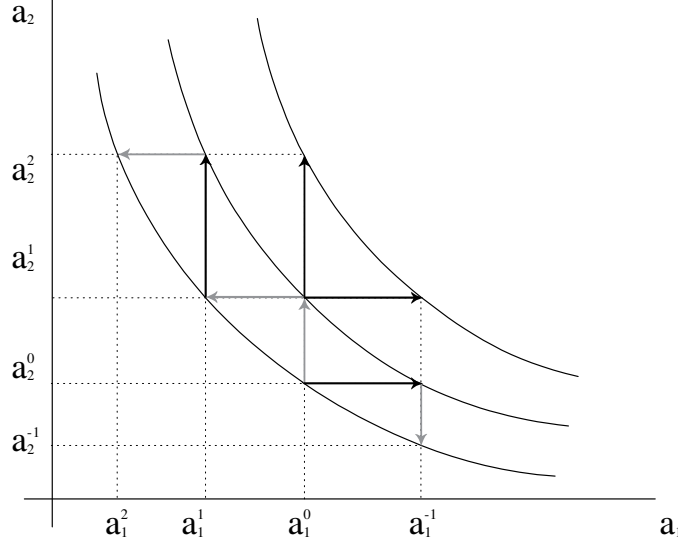


Figure 3: Lock-Step Procedure, Constructing a flight of stairs.

a_2 : $\left(a_1^{\frac{n+1}{2}}, a_2^{\frac{n+1}{2}}\right) \sim_f \left(a_1^{\frac{n}{2}}, a_2^{\frac{n}{2}}\right)$ and $\left(a_1^{\frac{n-1}{2}}, a_2^{\frac{n}{2}}\right) \sim_f \left(a_1^{\frac{n+1}{2}}, a_2^{\frac{n+2}{2}}\right)$. Therefore, by construction $\left(a_1^{\frac{n}{2}}, a_2^{\frac{n+1}{2}}\right) \sim_f \left(a_1^{\frac{n+1}{2}}, a_2^{\frac{n+2}{2}}\right)$ and then by F_4 , $\left(a_1^{\frac{n-1}{2}}, a_2^{\frac{n+1}{2}}\right) \sim_f \left(a_1^{\frac{n}{2}}, a_2^{\frac{n+2}{2}}\right)$.

Proceed analogously in the direction of increasing a_2 .

This demonstrates that if the vertical distance, measured in second component's units, between the indifference curves through (a_1^0, a_2^0) and (a_1^0, a_2^1) in a_1^n is the same as between those through (a_1^0, a_2^1) and (a_1^0, a_2^2) in a_1^{n-1} , then it is also the same between those through (a_1^0, a_2^0) and $(a_1^0, a_2^{\frac{1}{2}})$ in $a_1^{\frac{n}{2}}$ and between those through $(a_1^0, a_2^{\frac{1}{2}})$ and (a_1^0, a_2^1) in $a_1^{\frac{n-1}{2}}$. Repeating this procedure we can generate a dense set of points on indifference curves that are dense in \mathbb{R}_+^2 . Then continuity of \succ_f allows us to complete the entire map. Hence, if $(a_1, a_2) \sim_f (a'_1, a'_2)$ and $(a_1, \tilde{a}_2) \sim_f (a'_1, \tilde{a}'_2)$, then $(\tilde{a}_1, a_2) \sim_f (\tilde{a}'_1, a'_2) \Leftrightarrow (\tilde{a}_1, \tilde{a}_2) \sim_f (\tilde{a}'_1, \tilde{a}'_2)$.

As a result, we can create a mapping $a_2 \rightarrow \gamma(a_2)$ that transforms the original indifference map to be quasi-linear (vertically parallel indifference curves). The algorithm, which is formally described below, involves proceeding in infinitesimal steps and equalizing the step heights .

Set $\gamma(1) := 0$. To determine $\gamma(a_2)$ for $a_2 > 1$, pick an arbitrary a_1 and let a_1^0 solve $(a_1, a_2) \sim_f (a_1^0, 1 + \Delta)$, where Δ will be infinitesimal for the integration.³⁸ This solution exists by F_3 . Then for every $a_2^* \in (1, a_2]$,³⁹

Let a_1^* solve $(a_1^*, a_2^*) \sim_f (a_1^0, 1 + \Delta)$.

³⁸As established above, the result of this mapping will be independent of the choice of a_1 .

³⁹The existence of solutions in the two cases below is guaranteed by the same reasoning as in the above discussion.

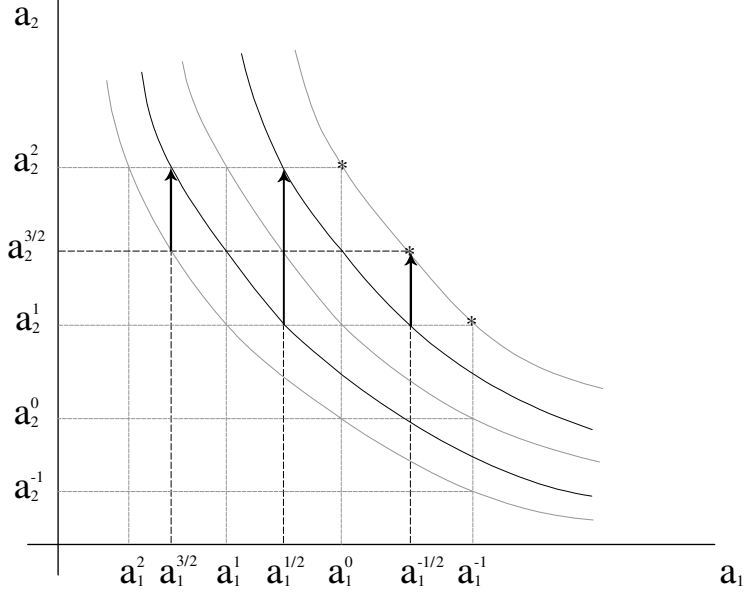


Figure 4: Lock-Step Procedure, Completing the indifference map.

Let a_1^{**} solve $(a_1^{**}, a_2^* + \Delta) \sim_f (a_1^0, 1 + \Delta)$.

Let a_2' solve $(a_1^*, a_2') \sim_f (a_1^0, 1)$.

Let da_2' solve $(a_1^{**}, a_2' + da_2') \sim_f (a_1^0, 1)$.

Note that by F_2 , $a_2' < a_2^*$ and $a_2' + da_2' < a_2^* + \Delta$.

Define implicitly $d\gamma(a_2^*) := \tilde{\gamma}(a_2' + da_2') - \gamma(a_2')$, where

$$\tilde{\gamma}(a) := \begin{cases} \gamma(a) & \text{for } a \leq a_2^* \\ \gamma(a_2^*) + a - a_2^* & \text{for } a > a_2^* \end{cases}$$

and then

$$\gamma(a_2) := \gamma(1) + \int_1^{a_2} d\gamma(a_2^*) = \int_1^{a_2} d\gamma(a_2^*).$$

Analogously determine $\gamma(a_2)$ for $a_2 < 1$: Pick an arbitrary a_1^0 and let a_1 solve $(a_1, a_2) \sim_f (a_1^0, 1)$. Then for every $a_2^* \in [a_2, 1)$:

Let a_1^* solve $(a_1^*, a_2^*) \sim_f (a_1^0, 1)$.

Let a_1^{**} solve $(a_1^{**}, a_2^* - \Delta) \sim_f (a_1^0, 1)$.

Let a_2' solve $(a_1^*, a_2') \sim_f (a_1^0, 1 + \Delta)$.

Let da_2' solve $(a_1^{**}, a_2' - da_2') \sim_f (a_1^0, 1 + \Delta)$.

Note that $a_2' < a_2^*$ and $a_2' + da_2' < a_2^* + \Delta$ by F_2 .

Define implicitly $d\gamma(a_2^*) := \gamma(a_2') - \tilde{\gamma}(a_2' - da_2')$, where

$$\tilde{\gamma}(a) := \begin{cases} \gamma(a) & \text{for } a \geq a_2^* \\ \gamma(a_2^*) - a + a_2^* & \text{for } a < a_2^* \end{cases}$$

and

$$\gamma(a_2) := \gamma(1) + \int_1^{a_2} d\gamma(a_2^*) = - \int_{a_2}^1 d\gamma(a_2^*) < 0.$$

Then $\gamma : \mathbb{R}_+ \rightarrow \mathbb{R}$, is a continuous and increasing function. The \succ_f indifference curves are quasi-linear with respect to $\gamma(a_2)$, so there is an increasing continuous function $\xi : \mathbb{R}_+ \rightarrow \mathbb{R}$, such that $\xi(a_1) + \gamma(a_2)$ generates the same indifference map. Hence re-defining

$$\varphi(a) := \xi(a_1) + \gamma(a_2)$$

represents \succ_f . Define

$$v_1(a_1) := \exp(\xi(a_1)) \text{ and } v_2(a_2) := \exp(\gamma(a_2)).$$

Then $v_1, v_2 : \mathbb{R}_+ \rightarrow \mathbb{R}_{++}$ are increasing and continuous and if we re-define, yet again, $\varphi(\mathbf{a}) := v_1(a_1)v_2(a_2)$, it represents \succ_f . By F_3 , the functions v_1, v_2 must be unbounded.

That the representation satisfies the axioms is easy to verify. ■

8.4. Proof of Theorem 4

(i) The analogue of Theorem 1 can be established by substituting \mathbf{a}_{-1} for a_2 in the theorem and in the proof, where now $\varphi : \mathbb{R}_+^N \rightarrow \mathbb{R}$.

Let φ be a representation of \succ_f . Let $\bar{\varphi} := \sup_{\mathbf{a} \in \mathbb{R}_+^N} \varphi(\mathbf{a})$ and $\underline{\varphi} := \inf_{\mathbf{a} \in \mathbb{R}_+^N} \varphi(\mathbf{a})$, if they are well defined. Otherwise, take $\bar{\varphi} = \infty$ and $\underline{\varphi} = -\infty$.

As before, let $S := \{(\mathbf{a}', \varphi') : Y(\mathbf{a}', \varphi') \neq \emptyset\}$. By F_3^N and the representation analogous to Theorem 1, $u(a_1) - u(0) > g(\mathbf{a}, \varphi)$ for $(\mathbf{a}, \varphi) \in S$.

Let \succ_S be a binary relation on S defined by $(\mathbf{a}, \varphi) \succ_S (\tilde{\mathbf{a}}, \tilde{\varphi}) \Leftrightarrow \{\mathbf{a}, \mathbf{b}\} \succ \{\tilde{\mathbf{a}}, \tilde{\mathbf{b}}\} \forall \mathbf{b} \in Y(\mathbf{a}, \varphi) \text{ and } \forall \tilde{\mathbf{b}} \in Y(\tilde{\mathbf{a}}, \tilde{\varphi})$.

Define $U_S : \mathbb{R}_+^N \times (\underline{\varphi}, \bar{\varphi}) \rightarrow \mathbb{R}$ such that $U_S(\mathbf{a}, \varphi) := U(\mathbf{a}, \mathbf{b})$ for some $\mathbf{b} \in Y(\mathbf{a}, \varphi)$. By Theorem 1, \succ_S is a weak order that can be represented by U_S . Note that the Consistency axiom (P_5) is relevant precisely on this domain. For $(\mathbf{a}, \varphi) \notin S$ define

$$U_S(\mathbf{a}, \varphi) := \begin{cases} 0 & \text{for } \varphi(\mathbf{a}) < \varphi \\ u(a_1) & \text{for } \varphi(\mathbf{a}) \geq \varphi \end{cases}.$$

Claim 4.1: U_S is continuous in all arguments.

Proof: Since the utility function is continuous on S , and because outside of S the function was chosen to be either a constant (hence continuous) or a continuous function, the only candidates for discontinuity are points on the boundary of S . There are two cases:

Case 1) $\varphi(\mathbf{a}) \geq \varphi$: Take $(\mathbf{a}, \varphi) \in bdr(S)$. Since (\mathbf{a}, φ) is a boundary point, it must be that $\varphi(\mathbf{a}) = \varphi$. Now let $\{\mathbf{a}^n, \varphi^n\}$ be a sequence in S which converges to (\mathbf{a}, φ) . By the definition of S , $U_s((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) = u(a_1^n) - g((a_1^n, \mathbf{a}_{-1}^n), \varphi^n)$. Because preferences are continuous and using the properties of g from Theorem 1, we have $\lim_{n \rightarrow \infty} u(a_1^n) - g((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) = u(a_1)$ as required.

Case 2) $\varphi(\mathbf{a}) < \varphi$: Take $(\mathbf{a}, \varphi) \in bdr(S)$. Again, let $\{\mathbf{a}^n, \varphi^n\}$ be an arbitrary sequence in S which converges to (\mathbf{a}, φ) . By the definition of S ,

$$U_s((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) = u(a_1^n) - g((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) > \inf_{\mathbf{b}} \{u(b_1) : \varphi(\mathbf{b}) = \varphi^n \text{ and } b_1 < a_1^n\}.$$

Since \succ is continuous, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} u(a_1^n) - g((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) &= u(a_1) - g((a_1, \mathbf{a}_{-1}), \varphi) \geq \\ \inf_{\mathbf{b}} \{u(b_1) : \varphi(\mathbf{b}) = \varphi \text{ and } b_1 < a_1\} &= u(0). \end{aligned}$$

where the last equality is implied by F_3^N . As $(\mathbf{a}, \varphi) \notin S$, we claim that $u(a_1) - g((a_1, \mathbf{a}_{-1}), \varphi) \leq \inf_{\mathbf{b}} \{u(b_1) : \varphi(\mathbf{b}) = \varphi \text{ and } \{\mathbf{b}\} \sim \{\mathbf{a}, \mathbf{b}\}\} = u(0)$. If not, then $u(a_1) - g((a_1, \mathbf{a}_{-1}), \varphi) = u(c_1) > u(0)$. But for any \mathbf{c} with $c_1 > 0$, using F_3^N , we could find \mathbf{c}' with $c'_1 < c_1$ and $\varphi(\mathbf{c}') = \varphi(\mathbf{c})$. Using Theorem 1, this would imply that $(\mathbf{a}, \varphi) \in S$, which is a contradiction. Combining we have $\lim_{n \rightarrow \infty} u(a_1^n) - g((a_1^n, \mathbf{a}_{-1}^n), \varphi^n) = u(0)$, as required. ||

Definition: For $(\mathbf{a}, \varphi) \in S$, define $I_S(\mathbf{a}, \varphi) := \{(\mathbf{a}', \varphi') : (\mathbf{a}', \varphi') \sim_S (\mathbf{a}, \varphi)\} \subseteq S$. That is, $I_S(\mathbf{a}, \varphi)$ is the \succ_S equivalence class of (\mathbf{a}, φ) .

Let $a_1^* : \mathbb{R}_+^2 \times (\underline{\varphi}, \bar{\varphi}) \rightarrow \mathbb{R}_+$ be the solution to

$$u(a_1^*(\mathbf{a}, \varphi)) = u(a_1) - g(\mathbf{a}, \varphi) = U_S(\mathbf{a}, \varphi).$$

a_1^* is the "first component equivalent" functional on S .⁴⁰ Since $u(a_1) > u(a_1) - g(\mathbf{a}, \varphi) >$

⁴⁰Formally, $\forall \mathbf{x} \in \mathbb{R}_+^{N-1}$, $\{(a_1^*(\mathbf{a}, \varphi), \mathbf{x})\} \sim \{\mathbf{a}, \mathbf{b}\}$, $\forall \mathbf{b} \in Y(\mathbf{a}, \varphi)$

$u(0)$ and \succ_S is continuous, a_1^* is well defined and we have $(\mathbf{a}, \varphi) \succ_S (\tilde{\mathbf{a}}, \tilde{\varphi}) \Leftrightarrow a_1^*(\mathbf{a}, \varphi) > a_1^*(\tilde{\mathbf{a}}, \tilde{\varphi})$.

Claim 4.2: The shame $g(\mathbf{a}, \varphi)$ is strictly increasing in φ .

Proof: Assume to the contrary that there is $\varphi' > \varphi$ and $(\mathbf{a}, \varphi') \sim_S (\mathbf{a}, \varphi)$ for some \mathbf{a} . Then for $\varphi' > \varphi'' > \varphi''' > \varphi$ we must have $(\mathbf{a}, \varphi'') \sim_S (\mathbf{a}, \varphi''')$ as shame is weakly increasing in φ . Now pick \mathbf{a}' such that $(\mathbf{a}', \varphi) \succ_S (\mathbf{a}', \varphi')$ and $(\mathbf{a}', \varphi), (\mathbf{a}', \varphi') \in S$. This is possible by continuity of U_S , since for \mathbf{a}'' such that $\varphi(\mathbf{a}'') = \varphi$ the definition of U_S yields $U_S(\mathbf{a}'', \varphi) > U_S(\mathbf{a}'', \varphi')$. Then by P_5 , $(\mathbf{a}', \varphi''') \succ_S (\mathbf{a}', \varphi'')$, a contradiction to shame being weakly increasing in φ .||

Claim 4.3: For all (\mathbf{a}, φ) and $\tilde{\varphi} \in (\varphi(a_1, \mathbf{0}), \bar{\varphi})$ there exists $\tilde{\mathbf{a}}$ such that $(\tilde{\mathbf{a}}, \tilde{\varphi}) \in I_S(\mathbf{a}, \varphi)$.

Proof: Define φ^* implicitly by $U_s((a_1, \mathbf{0}), \varphi^*) = U_s(\mathbf{a}, \varphi)$. This is possible by the Intermediate Value Theorem, as $U_s((a_1, \mathbf{0}), \varphi(a_1, \mathbf{0})) = u(a_1) > U_s(\mathbf{a}, \varphi) > U_s((a_1, \mathbf{0}), \varphi)$, where the last inequality is due to P_4 and Claim 4.2. There are two cases to consider:

Case 1) $\tilde{\varphi} \geq \varphi^*$: Then $U_s((a_1, \mathbf{0}), \tilde{\varphi}) \leq U_s(\mathbf{a}, \varphi)$ according to the monotonicity of shame. By F_3^N there is $\bar{a}_2(\tilde{\varphi})$ that solves $\varphi(a_1, \bar{a}_2(\tilde{\varphi}), \mathbf{0}) = \tilde{\varphi}$. Then $U_s((a_1, \bar{a}_2(\tilde{\varphi}), \mathbf{0}), \tilde{\varphi}) \geq U_s(\mathbf{a}, \varphi)$ and by the Intermediate Value Theorem there is $\tilde{a}_2(\tilde{\varphi}) \in [0, \bar{a}_2(\tilde{\varphi})]$ such that

$$U_s((a_1, \tilde{a}_2(\tilde{\varphi}), \mathbf{0}), \tilde{\varphi}) = U_s(\mathbf{a}, \varphi).$$

Case 2) $\tilde{\varphi} < \varphi^*$: Then

$$U_s((a_1^*(\mathbf{a}, \varphi), \mathbf{0}), \tilde{\varphi}) \leq U_s(\mathbf{a}, \varphi) \leq U_s((a_1, \mathbf{0}), \tilde{\varphi}).$$

By the Intermediate Value Theorem there is $\tilde{a}_1(\tilde{\varphi}) \in [a_1^*(\mathbf{a}, \varphi), a_1]$ such that

$$U_s((\tilde{a}_1(\tilde{\varphi}), \mathbf{0}), \tilde{\varphi}) = U_s(\mathbf{a}, \varphi).||$$

Combining the two cases we see that $\tilde{\varphi}$ parametrizes a path

$$\tilde{\mathbf{a}}_{(\mathbf{a}, \varphi)}(\tilde{\varphi}) := \begin{cases} (\tilde{a}_1(\tilde{\varphi}), \mathbf{0}) & \text{for } \tilde{\varphi} < \varphi^* \\ (a_1, \tilde{a}_2(\tilde{\varphi}), \mathbf{0}) & \text{for } \tilde{\varphi} \geq \varphi^* \end{cases}$$

of allocations. According to Claim 4.2 $\varphi(\mathbf{a})$ must be strictly increasing along this path. This

implies $\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ is strictly increasing in its first component for $\widetilde{\varphi} < \varphi^*$ and in its second component for $\widetilde{\varphi} \geq \varphi^*$.

Now we construct a \succ_S indifference class close to the original one:

Claim 4.4: For $\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ as defined above, $\widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ that solves

$$\left(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi}), \widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}) \right) \in I_S(\mathbf{a}, \varphi + d\varphi)$$

is increasing in $\widetilde{\varphi}$.

Proof: Assume $\widetilde{\varphi}' > \widetilde{\varphi}$. There are two cases to consider:

Case 1) $\widetilde{\varphi}' > \varphi^*$: Then $\widetilde{a}_{1(\mathbf{a},\varphi)}(\widetilde{\varphi}') = a_1$, $\widetilde{a}_{1(\mathbf{a},\varphi)}(\widetilde{\varphi}) \leq a_1$ and $\widetilde{a}_{2(\mathbf{a},\varphi)}(\widetilde{\varphi}') > \widetilde{a}_{2(\mathbf{a},\varphi)}(\widetilde{\varphi})$.

P_4 implies

$$\left(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi}), \widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}) \right) \prec_S \left(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi}'), \widetilde{\varphi}' + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}') \right).$$

Case 2) $\widetilde{\varphi}' \leq \varphi^*$: Then $\widetilde{a}_{2(\mathbf{a},\varphi)}(\widetilde{\varphi}') = \widetilde{a}_{2(\mathbf{a},\varphi)}(\widetilde{\varphi}) = 0$ and $\widetilde{a}_{1(\mathbf{a},\varphi)}(\widetilde{\varphi}') > \widetilde{a}_{1(\mathbf{a},\varphi)}(\widetilde{\varphi})$.

As \succ_S is increasing in a_1 ,

$$\left(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi}), \widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}) \right) \prec_S \left(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\widetilde{\varphi}'), \widetilde{\varphi}' + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}') \right).$$

As φ increases in φ , we must have $\widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi}') > \widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ in both cases. ||

Now we define a re-scaling $\varphi \mapsto \gamma(\varphi)$ in order to transform the original indifference map of $U_S(\mathbf{a}, \varphi)$ to be quasi-linear. We proceed similarly to the proof of Theorem 3. Choose $\varphi^0 \in (\underline{\varphi}, \overline{\varphi})$ and define $\gamma(\varphi^0) := 1$. Further set $\gamma(\varphi^0 + d\varphi) := 1 + d\gamma$, where $d\varphi$ is infinitesimal. To define $\gamma(\varphi)$ for $\varphi \neq \varphi^0$, pick \mathbf{a} such that $\varphi_{(\mathbf{a},\varphi)}^* < \varphi^0$. As $\varphi_{(\mathbf{a},\varphi)}^* < \varphi$, this implies $\varphi_{(\mathbf{a},\varphi)}^* < \min[\varphi, \varphi^0]$. Choose \mathbf{a}^0 such that $(\mathbf{a}^0, \varphi^0) \in I_S(\mathbf{a}, \varphi)$. We will look at the increasing graphs $\widetilde{\varphi}$ and $\widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ as defined above. Consider two cases for applying the Lock-Step Procedure:

Case 1) $\varphi > \varphi^0$: Define a climbing flight of stairs between the graphs $\widetilde{\varphi}$ and $\widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ recursively: Let φ^{n+1} solve $(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\varphi^n), \varphi^{n+1}) \sim_S (\mathbf{a}^0, \varphi^0 + d\varphi)$. The solution exists by the construction of $\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\varphi^n)$.

Case 2) $\varphi < \varphi^0$: Define a descending flight of stairs between the graphs $\widetilde{\varphi}$ and $\widetilde{\varphi} + d\varphi_{(\mathbf{a},\varphi)}(\widetilde{\varphi})$ recursively: Let φ^{-n-1} solve $(\widetilde{\mathbf{a}}_{(\mathbf{a},\varphi)}(\varphi^{-n-1}), \varphi^{-n}) \sim_S (\mathbf{a}^0, \varphi^0 + d\varphi)$.

Then $\gamma(\widetilde{\varphi})$ can be determined analogously to the proof of Theorem 2 by equalizing all step-heights to $d\varphi$ and integrating. Due to P_5 this definition is independent of the choice of \mathbf{a}^0 .

Now the indifference map of $U_S(\mathbf{a}, \varphi)$ is quasi linear in $\gamma(\varphi)$, where $\gamma : \mathbb{R}_{++} \rightarrow \mathbb{R}$ is strictly increasing and continuous. Further remember that $U_S(\mathbf{a}, \varphi)$ is strictly decreasing in φ . Therefore, there exists $H : \mathbb{R}_+^N \rightarrow \mathbb{R}$, such that $H(\mathbf{a}) - \gamma(\varphi)$ represents \succ_S on S .

Define $u_S(a_1) := H(\mathbf{a}) - \lim_{\varphi \rightarrow \varphi(\mathbf{a})} \gamma(\varphi)$. Because of P_1 ,

$$U(\{\mathbf{a}, \mathbf{b}\}) := \begin{cases} u_S(a_1) & \text{if } \{\mathbf{a}\} \sim \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\} \\ H(\mathbf{a}) - \gamma(\varphi(\mathbf{b})) & \text{if } \{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \succ \{\mathbf{b}\} \\ u_S(b_1) & \text{if } \{\mathbf{a}\} \succ \{\mathbf{a}, \mathbf{b}\} \sim \{\mathbf{b}\} \end{cases}$$

represents \succ confined to the collection of all two element sets. Therefore, $H(\mathbf{a}) \equiv u_S(a_1) + \gamma(\varphi(\mathbf{a}))$ must hold. Hence

$$U(A) = \max_{\mathbf{a} \in A} [u_S(a_1) + \gamma(\varphi(\mathbf{a}))] - \max_{\mathbf{b} \in A} [\gamma(\varphi(\mathbf{b}))]$$

represents \succ on K , where φ represents \succ_f , and u_s and γ are strictly increasing. Since φ represents \succ_f , so does $\gamma(\varphi)$. Hence, there is a representation φ of \succ_f , such that γ is the identity and

$$U(A) = \max_{\mathbf{a} \in A} [u_s(a_1) + \varphi(\mathbf{a})] - \max_{\mathbf{b} \in A} [\varphi(\mathbf{b})]$$

represents \succ on K .

(ii) To establish the analogue of Theorem 3, namely that there are N increasing unbounded functions v_1, \dots, v_N , such that the fairness ranking \succ_f can be represented by $\varphi(\mathbf{a}) = v_1(a_1) \cdot \dots \cdot v_N(a_N)$, if and only if it satisfies F_1, F_2, F_3^N and F_4^N we apply the Theorem of Luce and Tukey, just as in the proof of Theorem 3. It establishes the existence of an additive representation $\xi_1(a_1) + \dots + \xi_N(a_N)$ of \succ_f . Define $v_n(a_n) := \exp(\xi_n(a_n))$ for all $n \in \{1, \dots, N\}$. Then $v_1, \dots, v_N : \mathbb{R}_+ \rightarrow \mathbb{R}_{++}$ are increasing and continuous and if we re-define $\varphi(\mathbf{a}) := v_1(a_1) \cdot \dots \cdot v_N(a_N)$, it represents \succ_f . By F_3^N , the functions v_1, \dots, v_N must be unbounded.

That the representations satisfy the axioms is easy to verify. ■

8.5. Proof of Proposition

Define the random variable $X_n := \bar{\theta}_n - x_1$.

Claim 5.1: DM's utility is a decreasing function of $|X_n|$

Proof: DM chooses a to maximize

$$u(d(a, x_1)) + \beta_N \left(E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right] - E \left[\prod_{i=1}^N u(d(a^*, x_i)) \mid x^n \right] \right).$$

Since θ is a sufficient statistic for x^N , we can write:

$$E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right] = E \left[E \left[\prod_{i=1}^N u(d(a, x_i)) \mid \theta \right] \mid x^n \right].$$

This expression is single peaked as a function of a :

$$E \left[\prod_{i=1}^N u(d(a, x_i)) \mid \theta \right] = E \left[\prod_{i=1}^N u(d(a - \theta, x_i)) \mid 0 \right] =: f(a - \theta),$$

where f is symmetric and single peaked, with a peak in 0 and ⁴¹ Write $\pi_{x^n}(\theta)$ for the density function corresponding to $\theta \sim N(\bar{\theta}_n, \bar{\nu}_n^2)$. It is single peaked with peak in $\bar{\theta}_n$. Thus,

$$E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right] = \int_{\Omega} f(a - \theta) \pi_{x^n}(\theta) d\theta$$

is the convolution of two symmetric single peaked functions, with peak in 0 and $\bar{\theta}_n$, respectively. Then $E \left[\prod_{i=1}^N u(d(a, x_i)) \mid x^n \right]$ is single peaked with peak in $a = \bar{\theta}_n$. This means that fairness is maximized at $a^* = \bar{\theta}_n$ and, therefore, shame is increasing with $|\bar{\theta}_n - a|$. By assumption, DM's selfish utility is decreasing with $|a - x_1|$. Therefore, DM in effect chooses $|\bar{\theta}_n - a| \in [0, |X_n|]$. Fix X_n and denote by $l(|X_n|)$ DM's optimal choice of $|\bar{\theta}_n - a|$ and $V(l(|X_n|), |X_n|)$ the associated (total) utility. Then for $|X'_n| < |X_n|$,

$$\max \{0, |X'_n| - l(|X_n|)\} \leq |X_n| - l(|X_n|)$$

and

$$\min \{|X'_n|, l(|X_n|)\} \leq l(|X_n|)$$

with at least one inequality strict. Therefore,

$$V(l(|X_n|), |X_n|) < V(\min \{|X'_n|, l(|X_n|)\}, |X'_n|)$$

and by definition

$$V(\min \{|X'_n|, l(|X_n|)\}, |X'_n|) \leq V(l(|X'_n|), |X'_n|).$$

⁴¹The first equality is justified, since only the distance, which is symmetric, enters the utility functions.

Combining the two inequalities establishes the result.||

For given θ , note that

$$X_n | \theta \sim N \left(-\frac{\theta}{\nu^2 + n\sigma^2}, \frac{\sigma^2 (\sigma^2 (n-1) (2v^2 + n\sigma^2) + v^4)}{(n\sigma^2 + \nu^2)^2} \right).$$

Define $H_n(\theta) := -\frac{\nu^2}{\nu^2 + n\sigma^2}\theta$, so $H_n(\theta) \sim N \left(0, \frac{\nu^6}{(\nu^2 + n\sigma^2)^2} \right)$. Then

$$[X_n - H_n(\theta)] | H_n(\theta) \sim N \left(0, \frac{\sigma^2 (\sigma^2 (n-1) (2v^2 + n\sigma^2) + v^4)}{(n\sigma^2 + \nu^2)^2} \right).$$

Let $h_n(H_n(\theta))$ and $g_n(X_n - H_n(\theta))$ denote the associated density functions. The convolution

$$\int h_n(H_n(\theta)) g_n(X_n - H_n(\theta)) dH_n(\theta)$$

yields

$$X_n \sim N \left(0, \frac{\sigma^2 (\sigma^2 (m-1) (2v^2 + n\sigma^2) + v^4) + v^6}{(n\sigma^2 + \nu^2)^2} \right).$$

Hence for every n ,⁴²

$$p(|X_n| = z) \propto \exp \left[-\frac{z^2}{2} \frac{(\nu^2 + n\sigma^2)^2}{(\sigma^2 (\sigma^2 (n-1) (2v^2 + n\sigma^2) + v^4) + v^6)} \right].$$

Consequently

$$\begin{aligned} & \frac{p(|X_n| = z)}{p(|X_m| = z)} \\ & \propto \exp \left[-\frac{z^2}{2} \left(\frac{(\nu^2 + n\sigma^2)^2}{(\sigma^2 (\sigma^2 (n-1) (2v^2 + n\sigma^2) + v^4) + v^6)} - \frac{(\nu^2 + m\sigma^2)^2}{(\sigma^2 (\sigma^2 (m-1) (2v^2 + m\sigma^2) + v^4) + v^6)} \right) \right] \\ & = : \exp \left[-\frac{z^2}{2} c_{\nu, \sigma}(m, n) \right]. \end{aligned}$$

Thus, the ratio $\frac{p(|X_n|=z)}{p(|X_m|=z)}$ satisfies the Monotone Likelihood Ratio Property (MLRP): If $c_{\nu, \sigma}(m, n) < (>) 0$, then $\frac{p(|X_n|=z)}{p(|X_m|=z)}$ increases (decreases) with z . Since DM's utility is decreasing in z , she will strictly prefer m over n (n over m) if and only if $c_{\nu, \sigma}(m, n) < (>) 0$. In the text we define the ratio of the standard deviations σ and ν as $s := \frac{\sigma}{\nu}$. s is a suffi-

⁴²We use the fact that for a Normal distribution with mean 0, $p(X_n = z) = p(X_n = -z)$. With \propto we denote "proportional to".

cient statistic for (ν, σ) . Assuming $n > m$, we find, with some straightforward algebra, that $c_s(m, n) < (>) 0$, if and only if $2 + (m + n)s^2 - 3s^4 - 2(m + n)s^6 - mns^8 < (>) 0$.

Claim 5.2: For $n > m$, $c_s(m, m + 1) < 0$ implies $c_s(n, n + 1) < 0$.

Proof:

$$2 + (m + (m + 1))s^2 - 3s^4 - 2(m + (m + 1))s^6 - m(m + 1)s^8 < 0$$

is equivalent to

$$\frac{2 - 3s^4}{2m + 1} < 2s^6 + \frac{m^2 + m}{2m + 1}s^8 - s^2.$$

Let $lhs := \frac{2-3s^4}{2m+1}$ and $rhs := 2s^6 + \frac{m^2+m}{2m+1}s^8 - s^2$. Consider two cases:

- i) $2 - 3s^4 > 0$: $\frac{\partial}{\partial m}(lhs) < 0$ and $\frac{\partial}{\partial m}(rhs) > 0$ implies $c_s(n, n + 1) < 0$ for $n > m$.
- ii) $2 - 3s^4 \leq 0$: Then for all m , $lhs \leq 0$ and $rhs > 0$, which implies $c_s(n, n + 1) < 0$ for $n > m$.||

s^* as defined in the text solves $c_s(1, 2) = 0$.

Claim 5.3: For $s \leq s^*$, $c_s(m, m + 1) = 0$ has a unique solution $n^*(s) \in \mathbb{R}_+$. For $s > s^*$, no positive solution exists.

Proof: Assume first $s \leq s^*$: To show existence, note that due to the quadratic term in m , $m \rightarrow \infty$ implies $c_s(m, m + 1) \rightarrow -\infty$ for all s . Since $c_s(m, m + 1)$ is continuous in m , it is sufficient to show that $c_s(1, 2) > 0$ for all $s \leq s^*$. $c_s(1, 2) > 0$ is equivalent to

$$2 + 3s^2 - 3s^4 - 6s^6 - 2s^8 > 0$$

or

$$2/s^2 + 3 > 3s^2 + 6s^4 + 2s^6 \quad \forall s \leq s^*.$$

For this last inequality $\frac{\partial}{\partial s}(lhs) < 0$ and $\frac{\partial}{\partial s}(rhs) > 0$. Since $2/s^{*2} + 3s^* = 3s^{*2} + 6s^{*4} + 2s^{*6}$, $c_s(1, 2) > 0$ must hold for all $s \leq s^*$. Hence a solution $m = n^*(s)$ exists. Its uniqueness follows directly from claim 5.1.

Consider now the case $s > s^*$: In that case $c_s(1, 2) < 0$. Claim 1 implies that no solution exists to the equation $c_s(m, m + 1) = 0$ for any $m \in \mathbb{R}_+$.||

Claim 5.4: $n^*(s)$ is decreasing in s .

Proof:

$$2 + (n^* + (n^* + 1))s^2 - 3s^4 - 2(n^* + (n^* + 1))s^6 - n^*(n^* + 1)s^8 = 0$$

is equivalent to

$$\frac{2/s^2 - 3s^2}{2n^* + 1} = 2s^4 + \frac{n^{*2} + n^*}{2n^* + 1}s^6 - 1.$$

For this equality $\frac{\partial}{\partial n^*}(lhs) < 0$ and $\frac{\partial}{\partial n^*}(rhs) > 0$ for all s , while $\frac{\partial}{\partial s}(lhs) < 0$ and $\frac{\partial}{\partial s}(rhs) > 0$ for all n^* .||

Thus, if DM prefers $n = 1$ to $n = 2$, then $n^*(s) = 1$ is her globally preferred value. If she prefers m to both $m - 1$ and $m + 1$, then $n^*(s) = m$ is her globally preferred value. Neglecting the integer constraint, $m = n^*(s)$ is the unique positive solution to $2 + (m + (m + 1))s^2 - 3s^4 - 2(m + (m + 1))s^6 - m(m + 1)s^8 = 0$, if a solution exists. Furthermore, $n^*(s)$ is a decreasing function.■

References

- [1] Anderoni, James and John H. Miller (2002) "Giving according to GARP: An experimental test of the consistency of preference for altruism." *Econometrica*, 70, 737-753.
- [2] Benabou, Roland and Jean Tirole (2006) "Incentives and Prosocial Behavior." *American Economic Review*, 96(5), 1652-1678.
- [3] Bolton, Gary E., Elena Katok and Rami Zwick (1998) "Dictator game giving: Rules of fairness versus acts of kindness." *International Journal of Game Theory*, 27: 269-299.
- [4] Broberg, Thomas, Tore Ellingsen and Magnus Johannesson (2007) "Is Generosity Involuntary?" *Economics Letters*, 94, 32-37.
- [5] Burnham, Terence C. (2003) "Engineering altruism: A theoretical and experimental investigation of anonymity and gift giving." *Journal of Economic Behavior and Organization*, 50, 133-144.
- [6] Buss, Arnold H. (1980) "Self-Consciousness and Social Anxiety." San Francisco, W. H. Freeman.
- [7] Camerer, Colin (2003) "Behavioral Game Theory: Experiments in Strategic Interaction". Princeton University Press.

- [8] Charnes, Gary and Mathew Rabin (2002) "Understanding Social Preferences with Simple Tests." *Quarterly Journal of Economics*, 117 (3), 817-870.
- [9] Coate, Stephen and Stephen Morris (1995) "On the Form of Transfers to Special Interests." *Journal of Political Economy*, Vol. 103, No. 6, 1210-1235.
- [10] Dana, Jason D., Dalian M. Cain and Robin M. Dawes (2006) "What you don't Know Won't Hurt me: Costly (but quiet) Exit in a Dictator Game." *Organizational Behavior and Human Decision Processes*, 100(2), 193-201.
- [11] Dana, Jason D., Roberto A. Weber and Jason Xi Kuang (2005) "Exploiting moral wriggle room: Behavior inconsistent with a preference for fair outcomes." mimeo.
- [12] Davis, Douglas D. and Charles A. Holt (1993) "Experimental economics." Princeton University Press.
- [13] Dekel, Eddie, Barton L. Lipman and Aldo Rustichini (2005) "Temptation Driven Preferences." mimeo
- [14] Fehr Ernst, Klaus M. Schmidt (1999) "A Theory of Fairness, Competition and Cooperation." *Quarterly Journal of Economics*, 114, 817-868.
- [15] Frohlich, Norman, Joe Oppenheimer and J. Bernard Moore (2001) "Some doubts measuring self-interest using dictator experiments: The cost of anonymity." *Journal of Economic Behavior and Organization*, 46, 271-290.
- [16] Gauthier, David and Robert Sugden (editors) (June 1993) "Rationality, Justice and Social Contract: Themes from *Morals by Agreement*." University of Michigan Press.
- [17] Gul, Faruk and Wolfgang Pesendorfer (2005) "The Simple Theory of Temptation and Self-Control." mimeo
- [18] ——— (2001) "Temptation and Self Control." *Econometrica*, Vol. 69, No. 6, 1403-1435
- [19] Haley, Kevin J. and Daniel M.T. Fessler (2005) "Nobody's watching? Subtle cues affect generosity in an anonymous economic game." *Evolution and Human Behavior*, 26, 245-256.
- [20] Hammond, Peter J.(1991), "Interpersonal Comparisons of Utility: Why and How They Are and Should Be Made", in Elster and Roemer, "Interpersonal Comparisons of Well Being". Cambridge: Cambridge University Press, pp 200-254.

- [21] Johannesson, Magnus and Bjoran Persson (2000) "Non-reciprocal altruism in dictator games." *Economics Letters*, 69, 137-142.
- [22] Karni, Edi and Zvi Safra (1998) "The Hexagon Condition and Additive Representation for Two Dimensions: An Algebraic Approach." *Journal of Mathematical Psychology*, 42, 393-399.
- [23] Keeney, Ralph L. and Howard Raiffa, with a contribution by Richard F. Meyer (1976) "Decisions with multiple objectives : preferences and value trade-offs". New York : Wiley.
- [24] Koch, K. Alexander and Hans-Theo Norman (2005) "Giving in Dictator Games: Regard for Others or Regard by others?" mimeo.
- [25] Krantz, David H., R. Duncan Luce, Patrick C. Suppes and Amos Tversky (1971) "Foundations of Measurements, Vol 1." Academic Press, New York.
- [26] Lazear, Edward P., Ulrike Malmendier and Roberto A. Weber (2005) "Sorting in Experiments with Application to Social Preferences." mimeo.
- [27] Luce, R. Duncan and John W. Tukey (1964) "Simultaneous Conjoint Measurement: A New Type of Fundamental Measurement." *Journal of Mathematical Psychology*, 1,1-27.
- [28] Mariotti, Marco (1999) "Fair Bargains: Distributive Justice and Nash Bargaining Theory." *Review of Economic Studies*, Vol.66, 733-41.
- [29] Miller, Dale T. (1999) "The Norm of Self-Interest." *American Psychologist*, Vol. 54, No. 12, 1053-1060.
- [30] Nash, John F. (1953) "Two-Person Cooperative Games." *Econometrica*, Vol. 21, No. 1, 128-140.
- [31] — (1950) "The Bargaining Problem." *Econometrica*, Vol. 18, No. 2, 155-162.
- [32] Neilson, William S. (2006-a) "Axiomatic Reference Dependence in Behavior towards Others and Toward Risk." *Economic Theory* 28, 681-692.
- [33] — (2006-b) "A Theory of Kindness, Reluctance, and Shame in Dictator Games." mimeo.
- [34] Oberholzer-Gee, Felix and Reiner Eichenberger (2004) "Fairness in Extended Dictator Game Experiments". Working paper

- [35] Osborne, Martin J. and Ariel Rubinstein (1994) "A course in Game Theory." MIT press, ISBN 0-262-65040-1.
- [36] Pillutla, Madan M. and J. Keith Murningham (1995) "Being fair or appearing fair: Strategic behavior in ultimatum bargaining." *Academy of Management Journal*, 38,1408-1426.
- [37] Rawls, John (1971) "A Theory of Justice." *The Belknap Press of Harvard University Press*.
- [38] Tullock, Gordon (1983) "Economics of Income Redistribution." Boston: Kluwer-Nijhoff.
- [39] Wittman, Donald (1989) "Why Democracies Produce Efficient Results." *The Journal of Political Economy*, Vol. 97, No. 6, 1395-1424.