# Learning about a Class of Belief-Dependent Preferences without Information on Beliefs

Charles Bellemare

Alexander Sebald

Bellemare : Département d'économique, Université Laval, CIRPÉE
cbellemare@ecn.ulaval.ca
Sebald : Department of Economics, University of Copenhagen
alexander.sebald@econ.ku.dk

**Abstract:**
We show how to bound the effect of belief-dependent preferences on choices in sequential two-player games without information about the (higher-order) beliefs of players. The approach can be applied to a class of belief-dependent preferences which includes reciprocity (Dufwenberg and Kirchsteiger, 2004) and guilt aversion (Battigalli and Dufwenberg, 2007) as special cases. We show how the size of the bounds can be substantially reduced by exploiting a specific invariance property common to preferences in this class. We illustrate our approach by analyzing data from a large scale experiment conducted with a sample of participants randomly drawn from the Dutch population. We find that behavior of players in the experiment is consistent with significant guilt aversion: some groups of the population are willing the pay at least 0.16€ to avoid "letting down" another player by 1€. We also find that our approach produces narrow and thus very informative bounds on the effect of reciprocity in the games we consider. Our bounds suggest the model of reciprocity we consider is not a significant determinant of decisions in our experiment.

# 1   Introduction

In recent years there has been a growing interest in using belief-dependent preferences to explain experimental behavior ad odds with classical assumptions about human preferences (e.g. Charness and Dufwenberg (2006), Falk, Fehr, and Fischbacher (2008) and Charness and Dufwenberg (2010)). Belief-dependent preferences capture the idea that psychological factors such as people's beliefs concerning other people's intentions and expectations affect decision making.[1] Behavior may for example be motivated by the propensity to avoid feelings of guilt which result from 'letting down' others (see e.g. Battigalli and Dufwenberg, 2007). Guilt averse decision makers form beliefs about what others expect in order to infer how much these persons can be and are 'let down' by their own decisions. Alternatively, behavior may be motivated by reciprocity, i.e. the propensity to react kindly to perceived kindness and unkindly to perceived unkindness (see e.g. Dufwenberg and Kirchsteiger (2004)). Reciprocal decision makers form beliefs about the intentions of others in order to infer the (un)kindness of their behavior.

A natural approach to measure the relevance of belief-dependent preferences has been to test whether stated beliefs can predict behavior in a way consistent with a given type of belief-dependent preference. Charness and Dufwenberg (2006) for example ask players to state their higher-order beliefs in a trust game. They find that stated beliefs correlate with decisions in a way predicted by models of guilt aversion. More recently, Dhaene and Bouckaert (2010) measure the relevance of Dufwenberg and Kirchsteiger's (2004) theory of sequential reciprocity using stated first- and second-order beliefs and find empirical support.

Concerns have recently been expressed about the possibility that stated higher-order beliefs are correlated with preferences in a way which biases the estimated relevance of belief-dependent preferences. While beliefs and preferences may be correlated for various reasons, the source of this correlation is most often attributed to the presence of consensus effects which arise when individuals believe that others feel and think like themselves.

---

[1]Geanakoplos, Pearce, and Stacchetti (1989) and Battigalli and Dufwenberg (2009) present general frameworks to incorporate belief-dependent preferences in economics.

Consensus effects imply that stated beliefs about play in games correlate with preferences.[2] Vanberg (2010) uses a specific example to theoretically show that rational belief formation implies a correlation between preferences and beliefs of any order as long as preferences of players are correlated across the population. He concludes that we can expect such a correlation in experimental settings even when behavior is not driven by belief-dependent preferences. Bellemare, Sebald, and Strobel (2011) empirically investigate how correlation between preferences and stated beliefs can affect the estimated willingness to pay to avoid feeling guilty of letting down another player. They estimate this correlation by jointly modeling decisions and beliefs of players in a sequential trust game. They find that correlation between preferences and stated beliefs can exaggerate the measured level of guilt aversion in a population by a factor of two. Blanco, Engelmann, Koch, and Normann (2011) analyze the interaction between preferences and beliefs in a sequential prisoner's dilemma. They exploit data from a within-subject design (with participants playing both roles) and vary the information provided to players about the play of others to separately identify the direct impact of beliefs on decisions from consensus effects. They conclude that consensus effects are the primary determinants of the observed correlation between stated beliefs and decisions. These results highlight the complexity of measuring the relevance of belief-dependent preferences when exploiting data on higher-order beliefs.

This paper presents a new approach to learn about the relevance of belief-dependent preferences which does not require information about beliefs of players. We formally characterize conditions under which our approach can be used and we propose a simple two step estimation procedure to perform the required inferences. Although our approach is not exclusively tailored for experimental investigations, the conditions which have to be satisfied in order to use the approach proposed make it more suitable for controlled environments. Hence, we illustrate our approach by conducting an experiment using simple binary sequential two-player games to analyze the relevance of belief-dependent guilt aversion and reciprocity.

---

[2]Charness and Dufwenberg (2006) discuss the possibility that false consensus effects explain the correlation between decisions and beliefs in their data.

Our approach builds on random utility models to interpret the decisions of players in games.[3] We specify the utility of players as a function of their own monetary payoffs, their psychological payoffs which capture their belief-dependent preferences, as well as other unobservable factors. Our main parameter of interest measures the effect of belief-dependent preferences on behavior. Importantly, the psychological payoffs capturing the belief-dependent preferences are unknown variables without exploiting information on the beliefs of players. However, they are known to lie within well defined intervals. Our empirical strategy is to determine what can be learned about belief-dependent preferences from observing the monetary payoffs and the intervals of the psychological payoffs.

An immediate consequence of interval-measurements of the psychological payoffs is that the model parameters are set rather than point identified (see Manski and Tamer (2002)). Set identification implies that a range of parameter values – the identification region – are consistent with the data given the assumed model. The informativeness of the data given the model naturally decreases with the size of the identification region. Existing work has established that identification regions of the parameters of random utility models with interval measured regressors can be large and uninformative. Manski (2010) theoretically analyzes the binary random expected utility model when researchers do not have any information about the expectations of decision makers. He finds that the identification region of the model parameters is unbounded and thus uninformative when researchers cannot a priori sign the difference in expectation across both choices. Bellemare, Bissonnette, and Kröger (2010) analyze empirically decisions of senders in a binary trust game and estimate largely uninformative identification regions of their parameters when they do impose a priori assumptions about the beliefs of players. These results reveal the important difficulties confronting researchers interested in making inferences on belief-dependent preferences without using information about the beliefs of players.

One of the main insights of our analysis is that several prominent belief-dependent preferences satisfy an 'invariance property' which can be exploited to substantially reduce

---

[3]Random utility models have been extensively used to analyze choice behavior in experiments. See Cappelan, Hole, Sørensen, and Tungodden (2007), Bellemare, Kröger, and van Soest (2008).

the size of the identification region in order to produce informative bounds on the relevance of belief-dependent preferences. This invariance property is best described in the context of a game with two players – $A$ and $B$. The invariance property holds if player $B$'s decision is unaffected by his/her belief-dependent preferences when his choice cannot influence the final payoff of player $A$. To illustrate, suppose player $B$ must choose between two final allocations, both of which provide player $A$ with the same material payoff. Then, prominent models of guilt aversion predict that player $B$ cannot feel any guilt from letting down player $A$ by choosing a specific allocation because player $A$'s final payoff is independent of player $B$'s choice. Similarly, player $B$ cannot act reciprocally if player $A$'s payoff is independent of player $B$'s choice. This is because player $B$ cannot be (un)kind by providing player $A$ with an (below) above average payoff. It follows that our empirical strategy involves implementing a sufficiently high number of games in which the invariance property holds to identify and estimate all other model parameters but the sensitivity parameter measuring the relevance of belief-dependent preferences. We show that the ability to recover separate estimates of the remaining parameters is the key to reduce the size of the identification region and to obtain informative bounds on the relevance of belief-dependent preferences.

We illustrate our approach by conducting an experiment using simple binary sequential two-player games. We derive closed form expressions for the bounds of the sensitivity parameters measuring the relevance of simple guilt aversion (Battigalli and Dufwenberg (2007)) and reciprocity (Dufwenberg and Kirchsteiger (2004)) in this binary choice setting. We implement our experiment using the LISS panel, a large-scale Internet panel whose respondents form a representative sample of the Dutch population.[4] Close to 1500 panel members completed our experiment which involved 500 payoff-wise unique games. One third of these games satisfied the payoff invariance condition discussed above. We exploit the unique features of the panel to perform inferences for different socio-economic

---

[4]Other experimental studies which have used similar platforms to obtain a representative sample of participants from the Dutch population include Bellemare and Kröger (2007), Bellemare, Kröger, and van Soest (2008).

groups, allowing us to asses the heterogeneity in belief-dependent preferences across a broad population.

Our analysis of guilt aversion suggests that the population willingness to pay to avoid letting down the other player by 1€ is significantly different from zero and at least greater or equal to 0.08€. We also find that the lower bound of the willingness to pay to avoid guilt is higher for several groups of the population. In particular, we find that high educated individuals are willing to pay at least 0.14€ to avoid letting down the other player by 1€, while men are willing to pay at least 0.16€ to avoid letting down the other player by 1€. Our approach also produces very narrow and thus highly informative bounds around the relevance of reciprocity in our experiment. Our results suggest that reciprocity weakly predicts the final decisions made in our experiment for all groups of the population we consider. The narrowness of these bounds also suggests that stated belief data are not needed to make precise inferences on the relevance of reciprocity in our experiment.

The organization of the paper is as follows. Section 2 presents a class of two-player extensive form games with belief-dependent preferences (also called 'psychological games') which is used to formally characterize the conditions under which our approach can be used. Section 3 presents our proposed approach and details how it can be applied to the analysis of guilt aversion and reciprocity. Section 4 describes our experiment, data, and presents results for the analysis of guilt aversion and reciprocity. Section 5 discusses the possibility to make inferences at the individual level and concludes.

# 2  A class of psychological games

In this section we present a class of two-player extensive form games with belief-dependent preferences building on Battigalli and Dufwenberg (2009). This class of games represents the strategic environment and class of preferences which we use in the subsequent sections to formally characterize the conditions under which our identification approach can be used.

Formally, let the set of players be $\mathcal{N} = \{1, 2\}$. Denote as $\mathcal{H}$ the finite set of histories

$h$, with the empty sequence $h^0 \in \mathcal{H}$, and $Z$ the set of terminal histories. Histories $h \in \mathcal{H}$ are sequences that describe the choices of players on the path to history $h$. More precisely, a history of length $x \in X$ is a sequence of actions $h = (a^1, \ldots, a^x)$ where each $a^t = (a_1^t, a_2^t)$ represents the profile of actions taken at stage $t$ $(1 \leq t \leq x)$. The history $\tilde{h} = (\tilde{a}^1, \ldots, \tilde{a}^v)$ precedes $h = (a^1, \ldots, a^x)$, written $\tilde{h} < h$, if $\tilde{h}$ is a prefix of $h$ (i.e. $v < x$ and $(\tilde{a}^1, \ldots, \tilde{a}^v) = (a^1, \ldots, a^v)$). Turned up side down, we say that $h$ *immediately succeeds* $\tilde{h}$ if $x = v + 1$ and $(\tilde{a}^1, \ldots, \tilde{a}^v) = (a^1, \ldots, a^v)$. At each non-terminal history $h \in \mathcal{H} \backslash Z$ each player $i \in \mathcal{N}$ has a nonempty, finite set of feasible actions $\mathcal{A}_{i,h}$. A typical element of $\mathcal{A}_{i,h}$ is denoted by $a_{i,h}$. Note, $\mathcal{A}_{i,h}$ can be a singleton, meaning that player $i$ is inactive at history $h$. In fact, we assume that players do not choose simultaneously. Whenever a player $i \in \mathcal{N}$ is active, player $j \in \mathcal{N}$ with $j \neq i$ is inactive. Let $\mathcal{A}_i$ be the finite set of pure strategies of player $i \in \mathcal{N}$ and $\mathcal{A} = \prod_{i \in \mathcal{N}} \mathcal{A}_i$ be the set of joint pure strategies. Pure strategies of player $i \in \mathcal{N}$ and joint pure strategies are respectively denoted by $a_i = (a_{i,h})_{h \in \mathcal{H} \backslash Z}$ and $a$. Furthermore, denote by $\mathcal{A}_i(h)$ the set of strategies $a_i$ of player $i$ that allow for history $h$. Terminal histories $\zeta(a) \in Z$ depend on joint pure strategies $a \in \mathcal{A}$.

To capture belief-dependent preferences like reciprocity and guilt we assume that in every history players hold

(i) a belief about the strategies of the other player,

(ii) a belief about the belief of the other player,

(iii) a belief about the belief about the belief of the other player etc, and

(iv) we assume players update their beliefs as events unfold.

More specifically, we assume that players hold infinite hierarchies of conditional beliefs. This means, player $i \in \mathcal{N}$ holds an updated and revised belief $\mu_i^1(\cdot|h) \in \Delta(A_j(h))$ about the strategies of the co-player $j$, where $\Delta(A_j(h))$ denotes the set of behavioral strategies defined on $A_j(h)$. Given this, $\mu_i^1 = (\mu_i^1(\cdot|h))_{h \in \mathcal{H}}$ represents the system of first-order beliefs of player $i$. In addition, at every history $h$ player $i$ forms expectations $\mu_i^2(h) \equiv \mathbf{E}_i(\mu_j^1)$ over

the system of first-order beliefs of player $j$, forms expectations $\mu_i^3(h)$ over the system of second-order beliefs of player $j$ etc. $\mu_i^2$ and $\mu_i^3$ respectively denote the system of second- and third-order beliefs of player $i$. More generally, we denote the infinite hierarchy of conditional probability systems of player $i \in \mathcal{N}$ by $\boldsymbol{\mu}_i$ where $\boldsymbol{\mu}_i = (\mu_i^k)_{k=1}^\infty$ and a hierarchy of conditional probability systems up to order $k < \infty$ as $\boldsymbol{\mu}_i^k$.

Given this, we can define belief-dependent preferences:[5]

**Definition 1** *The belief-dependent utility $u$ of any player $i \in \mathcal{N}$ from choosing strategy $a_i \in \mathcal{A}_i(h)$ in history $h$ is:*

$$u_i(a_i, \boldsymbol{\mu}_i^k) = \pi_i(a_i, \mu_i^1) + \phi B_i(a_i, \boldsymbol{\mu}_i^k) + \lambda \epsilon_i(a_i)$$

*where $i, j \in \mathcal{N}$.*

First, $\pi_i(a_i, \mu_i^1(h))$ denotes player $i$'s expected material payoff depending on his strategy and first-order belief. Second, $B_i(a_i, \boldsymbol{\mu}_i^k)$ denotes player $i$'s belief-dependent psychological payoff depending on his strategy and his hierarchy of beliefs $\boldsymbol{\mu}_i^k$ up to order $k < \infty$. The belief-dependent payoff $B_i(a_i, \boldsymbol{\mu}_i^k)$ can, for example, capture belief-dependent reciprocity as defined by Dufwenberg and Kirchsteiger (2004) or simple guilt aversion as defined by Battigalli and Dufwenberg (2007). Our parameter of interest $\phi$ captures player $i$'s sensitivity to his/her psychological payoff. Lastly, $\epsilon_i(a_i)$ denotes unobserved preferences from choosing the strategy $a_i$ assumed to be independent of all variables entering the model, and $\lambda$ denotes a noise parameter.

This brings us to the definition of our class of psychological games

**Definition 2** *A two-player extensive form game with belief-dependent preferences is a tuple $\Gamma = \langle N, \mathcal{H}, (u_i)_{i \in N} \rangle$ with $u_i$ as defined in Definition 1.*

We next present our approach to make inferences on $\phi$ without information on $\boldsymbol{\mu}_i^k$.

---

[5]Note, in accordance with the existing literature on belief-dependent preferences we assume that overall utilities of players are additive in own material and belief-dependent psychological payoffs. In the framework of Battigalli and Dufwenberg (2009) a more general specification is presented in their definition 4 [p. 12].

# 3 The proposed approach

Our objective is to understand whether and under what conditions one can make inferences on the sensitivity parameter $\phi$ without information about the set of beliefs $\boldsymbol{\mu}_i^k$. Consider first the following condition in the context of the class of extensive form games with belief-dependent preferences $\Gamma$ defined in the previous section:

**History condition (H)** There exists a non-empty set $\tilde{\mathcal{H}} \subseteq \mathcal{H}$ such that

$$\tilde{\mathcal{H}} = \{h \in \mathcal{H} \ : \ A_{i,h} = \{a'_{i,h}, a''_{i,h}\} \text{ and } h \text{ is } \textit{immediately succeeded}$$
$$\text{by only terminal histories } \zeta(a) \in Z\}.$$

An extensive form game satisfies condition H if there exists a non-empty set of *last* non-terminal histories $\tilde{\mathcal{H}} \subseteq \mathcal{H}$ in which the active player has two pure actions to choose from. Consider, for example, the two player sequential game in Figure 1.

[Figure 1]

In this simple two-player game, player $j$ can choose either the outside option $R$ and determine the monetary payoffs for himself and player $i$, or he can choose $L$ and let player $i$ decide the final allocation. If player $j$ chooses the outside option $R$, then $\pi_i(R)$, $\pi_j(R)$ respectively denote the monetary payoffs of players $i$ and $j$. On the other hand, if he chooses $L$, then player $i$ must choose between $l$ and $r$ at history $h^1$ in the game tree. We denote as $\pi_z(r)$ and $\pi_z(l)$ for $z \in \{i, j\}$ the monetary payoffs when playing $r$ and $l$ at history $h^1$. Both histories $h^0$ and $h^1$ are histories in which the respective active player can choose between two pure actions, but only history $h^1$ is immediately succeeded by only terminal histories, i.e. $\tilde{\mathcal{H}} = \{h^1\}$. We say a game does not satisfy condition H whenever there is no last non-terminal history in this game in which the active player in that history only has two actions.

Let condition H hold. It follows that the utility of any player $i \in \mathcal{N}$ in history $\tilde{h} \in \tilde{\mathcal{H}}$ reduces to

$$u_i(a_{i,\tilde{h}}, \boldsymbol{\mu}_i^k) = \pi_i(a_{i,\tilde{h}}) + \phi B_i(a_{i,\tilde{h}}, \boldsymbol{\mu}_i^k) + \lambda \epsilon_i(a_{i,\tilde{h}})$$

where $a_{i,\tilde{h}} \in \{a'_{i,\tilde{h}}, a''_{i,\tilde{h}}\}$ and $k < \infty$. Note that $\pi_i(a_i, \mu_i^1) = \pi_i(a_i)$ since player $i$ no longer faces uncertainty over decisions of player $j$. Given this, define $\Delta u_{i,\tilde{h}} \equiv u_i(a'_{i,\tilde{h}}) - u_i(a''_{i,\tilde{h}})$, $\Delta\pi_{i,\tilde{h}} = \pi_i(a'_{i,\tilde{h}}) - \pi_i(a''_{i,\tilde{h}})$, and $\Delta\epsilon_{i,\tilde{h}} = \epsilon_i(a'_{i,\tilde{h}}) - \epsilon_i(a''_{i,\tilde{h}})$. Assuming expected utility maximization, player $i$ will choose to play $a'_{i,\tilde{h}}$ in history $\tilde{h}$ if

$$\Delta u_{i,\tilde{h}} = \Delta\pi_{i,\tilde{h}} + \phi\Delta B_{i,\tilde{h}} + \lambda\Delta\epsilon_{i,\tilde{h}} > 0 \tag{1}$$

The decision rule (1) leads to the following choice probability

$$\Pr(c_i = a'_{i,\tilde{h}}|\boldsymbol{\pi}, \boldsymbol{\mu}_i^k) = F\left([\Delta\pi_{i,\tilde{h}} + \phi\Delta B_{i,\tilde{h}}]/\lambda\right) \tag{2}$$

where $\boldsymbol{\pi}$ is a vector of payoffs for all $a_{i,\tilde{h}} \in A_{i,\tilde{h}}$ and $F(\cdot)$ denotes the cumulative distribution function of $\Delta\epsilon_{i,\tilde{h}}$. This is a standard binary choice model when beliefs $\Delta B_{i,\tilde{h}}$ are observed for all $i$. In the later case, estimation of the model parameters can be performed by assuming a specific parametric distribution for $F(\cdot)$ (eg. normal of logistic). Alternatively, semiparametric estimation of the parameters is possible (up to some normalization) by treating $F(\cdot)$ as an unknown nonparametric function (see eg. Klein and Spady (1993)). Unfortunately, the lack of information on $\boldsymbol{\mu}_i^k$ implies that $\Delta B_{i,\tilde{h}}$ is not observed. Hence, conventional parametric and semiparametric estimators of binary choice models cannot be used to make inferences on $\phi$. Define $\underline{\Delta B_{i,\tilde{h}}} = \inf_{\boldsymbol{\mu}_i^k} \Delta B_{i,\tilde{h}}$ and $\overline{\Delta B_{i,\tilde{h}}} = \sup_{\boldsymbol{\mu}_i^k} \Delta B_{i,\tilde{h}}$. It follows that without information on $\boldsymbol{\mu}_i^k$,

$$\Delta B_{i,\tilde{h}} \in [\underline{\Delta B_{i,\tilde{h}}}, \overline{\Delta B_{i,\tilde{h}}}] \tag{3}$$

Consider the case where $\phi \geq 0$. Then, it follows from equation 3 and the proof of Proposition 4 in Manski and Tamer (2002) that the following must hold for all $i$

$$\Pr(c_i = a'_{i,\tilde{h}}|\boldsymbol{\pi}, \underline{\Delta B_{i,\tilde{h}}}, \overline{\Delta B_{i,\tilde{h}}}) \in [F\left([\Delta\pi_{i,\tilde{h}} + \phi\underline{\Delta B_{i,\tilde{h}}}]/\lambda\right), F\left(\Delta\pi_{i,\tilde{h}} + \phi\overline{\Delta B_{i,\tilde{h}}}]/\lambda\right)] \tag{4}$$

Inverting $\Pr(c_i = a'_{i,\tilde{h}}|\boldsymbol{\pi}, \underline{\Delta B_{i,\tilde{h}}}, \overline{\Delta B_{i,\tilde{h}}})$ in (4) yields an equivalent and useful expression given by

$$\Delta\pi_{i,\tilde{h}} + \phi\underline{\Delta B_{i,\tilde{h}}} \leq Q_i\lambda \leq \Delta\pi_{i,\tilde{h}} + \phi\overline{\Delta B_{i,\tilde{h}}} \tag{5}$$

9

where $Q_i \equiv F^{-1}(\Pr(c_i = a'_{i,\tilde{h}} | \boldsymbol{\pi}, \underline{\Delta B_{i,\tilde{h}}}, \overline{\Delta B_{i,\tilde{h}}}))$. The identification region consists of all values $(\phi, \lambda)$ which are consistent with either (4) or (5) for all $i$. Our particular focus is on the identification region for $\phi$. The bounds in (5) fall in the class of monotone-index models with interval regressors analyzed in Manski and Tamer (2002). They have established in the Corollary to their Proposition 4 that the identification region for $(\phi, \lambda)$ is convex. Our particular focus in this paper is on the identification region of $\phi$. Manski (2010) analyzes the identification in a binary choice monotone-index model when one covariate is not observable due to lack of information on expectations.

There is a range of values of $\phi$ which satisfy (5) for each game. The identification region is given by the intersection of the ranges across all games. To characterize our main result, consider a set of games $\boldsymbol{\Gamma}$ with each $\Gamma \in \boldsymbol{\Gamma}$ satisfying condition H, a history $\tilde{h} \in \tilde{\mathcal{H}}$ and define the following 5 mutually exclusive dummy variables distinguishing the 5 types of games possibly present in the set.

$$
\begin{aligned}
d_i^1 &= \mathbf{1}(\underline{\Delta B_{i,\tilde{h}}} > 0, \overline{\Delta B_{i,\tilde{h}}} > 0) \\
d_i^2 &= \mathbf{1}(\underline{\Delta B_{i,\tilde{h}}} < 0, \overline{\Delta B_{i,\tilde{h}}} < 0) \\
d_i^3 &= \mathbf{1}(\underline{\Delta B_{i,\tilde{h}}} < 0, \overline{\Delta B_{i,\tilde{h}}} > 0) \\
d_i^4 &= \mathbf{1}(\underline{\Delta B_{i,\tilde{h}}} = 0, \overline{\Delta B_{i,\tilde{h}}} > 0) \\
d_i^5 &= \mathbf{1}(\underline{\Delta B_{i,\tilde{h}}} < 0, \overline{\Delta B_{i,\tilde{h}}} = 0)
\end{aligned}
\tag{6}
$$

such that $\sum_{j=1}^{5} d_i^j = 1$ for all $i$ and where $\mathbf{1}(A)$ denotes the indicator function taking a value of 1 when event $A$ occurs, and 0 otherwise. Let $Q_i = F^{-1}(\Pr(c_i = a'_i | \boldsymbol{\pi}_i))$ and

$$
\begin{aligned}
\phi_i^A &= \left( Q_i \lambda - \Delta \pi_{i,\tilde{h}} \right) / \overline{\Delta B_{i,\tilde{h}}} \\
\phi_i^B &= \left( Q_i \lambda - \Delta \pi_{i,\tilde{h}} \right) / \underline{\Delta B_{i,\tilde{h}}} \\
\phi_i^C &= \max \left\{ \phi_i^A, \phi_i^B \right\}
\end{aligned}
\begin{aligned}
&\tag{7} \\
&\tag{8}
\end{aligned}
$$

Given this we can state our main proposition:

**Proposition 1** *Consider a set of games $\boldsymbol{\Gamma}$ such that condition H is satisfied and a history $\tilde{h} \in \tilde{\mathcal{H}}$. Assume $\phi \geq 0$ and let $[\phi_\lambda^l, \phi_\lambda^u]$ denote the identification region of $\phi$ conditional on*

$\lambda$. Furthermore, let $\mathcal{D} \subset \mathbf{\Gamma}$ denote the set containing games with $d_i^1 = 1$ and games with $d_i^2 = 1$.

Then, the endpoints of the identification region are given by:[6]

$$\phi_\lambda^l = \max_{\forall i}[\max[\underline{\phi}_i, 0]] \tag{9}$$

$$\phi_\lambda^u = \min_{i \in \mathcal{D}}[\max[\overline{\phi}_i, 0]] \text{ if } \mathcal{D} \text{ is not empty} \tag{10}$$

$$= +\infty \text{ otherwise}$$

where

$$\underline{\phi}_i = \left(d_i^1 + d_i^4\right)\phi_i^A + \left(d_i^2 + d_i^5\right)\phi_i^B + d_i^3\phi_i^C$$

$$\overline{\phi}_i = d_i^1\phi_i^B + d_i^2\phi_i^A.$$

**Notes**. This proposition reveals that the identification region is given by the intersection of $[\underline{\phi}_i, \overline{\phi}_i]$ across all games, where $\underline{\phi}_i$ and $\overline{\phi}_i$ denote the lowest and highest values of $\phi$ consistent with the game played by player $i$ conditional on $\lambda$. Which of $\phi_i^A$, $\phi_i^B$, and $\phi_i^C$ will be used to compute $\underline{\phi}_i$ and $\overline{\phi}_i$ will depend on the signs of $\underline{\Delta B_{i,\tilde{h}}}$ and $\overline{\Delta B_{i,\tilde{h}}}$. Take games with $d_i^1 = 1$ and let $\phi \to 0$. It follows that the upper bound in (5) will equate $Q_i\lambda$ when $\phi = \phi_i^A$. This determines the lowest value of $\phi$ consistent with that game. Now let $\phi \to \infty$. It follows that the lower bound in (5) will equate $Q_i\lambda$ when $\phi = \phi_i^B$. This determines the highest value of $\phi$ consistent with that game. A similar analysis applies to the other four game types. We also note that $\max[\underline{\phi}_i, 0]$ and $\max[\overline{\phi}_i, 0]$ enter (9) and (10) to enforce the restriction that $\phi \geq 0$.

It follows from the proposition that knowledge of $\lambda$ can reduce substantially the identification region and thus allows for more precise inferences on the relevance of belief-dependent preferences. The main insight of the paper is summarized in the following three conditions.

**Invariance condition (I)** $\Delta B_i(\boldsymbol{\mu}_i^k) = 0$ when $\pi_j(a'_{i,\tilde{h}}) = \pi_j(a''_{i,\tilde{h}})$.

---

[6]The proposition is stated for $\phi \geq 0$. The case of $\phi \leq 0$ follows analogously with the endpoints of the identification region given by $\phi_\lambda^l = \max_{i \in \mathcal{D}}[\min[\underline{\phi}_i, 0]]$ if $\mathcal{D}$ is not empty and $\phi_\lambda^l = -\infty$ otherwise, while $\phi_\lambda^u = \min_{\forall i}[\min[\overline{\phi}_i, 0]]$.

11

**Support condition (S)** $\Pr(\pi_j(a'_{i,\tilde{h}}) = \pi_j(a''_{i,\tilde{h}})) > 0$.

**Noise condition (N)** $\lambda$ is independent of $\boldsymbol{\pi}$.

Condition I states that the difference between the psychological payoffs of player $i$ from his two actions $a'_{i,\tilde{h}}$ and $a''_{i,\tilde{h}}$ is zero if the payoffs of player $j$ do not vary with the action chosen by player $i$. This condition holds for several important preferences discussed in the literature (see sections 3.1 and 3.2 below). Condition S states that games where condition I holds should be present in the data. Note that such games can easily be implemented in an experiment by appropriate selection of player $j$ payoffs. Condition N states that the noise parameter does not vary with the payoffs of the game. It can however depend on the observable characteristics of players. Condition N implies that the value of $\lambda$ for games which satisfy condition S is the same as the corresponding noise level present in games with some payoff variation for player $j$. Supportive evidence for condition N can be obtained by estimating a reduced form version of equation (2), allowing the noise parameter to vary with the payoffs levels. Section 4.1 discusses this in more detail.

Together, conditions I, S and N allow separate identification of $\lambda$. In particular, for preferences satisfying condition I in some history $\tilde{h} \in \tilde{\mathcal{H}}$, it follows from (2) that the choice probabilities for games satisfying condition S are given by

$$\Pr(c_i = a'_{i,\tilde{h}} | \boldsymbol{\pi}_i) = F(\overline{\pi}_i / \lambda) \tag{11}$$

where the psychological payoffs drop out of the choice probabilities when $\pi_j(a'_{i,\tilde{h}}) = \pi_j(a''_{i,\tilde{h}})$. Equation (11) can thus be used to estimate $\lambda$ using only games which satisfy condition S.

We thus propose a simple two step estimation procedure. In the first step, we estimate $\lambda$ for preferences satisfying condition I using data from games satisfying condition S. In the second step we estimate the identification region $[\phi^l_{\hat{\lambda}}, \phi^u_{\hat{\lambda}}]$, conditional on the first step estimate of $\lambda$. We next discuss in detail prominent examples of belief dependent preferences which can be tested using our proposed two step procedure.

## 3.1 Example 1: guilt aversion ($\phi \leq 0$)

Battigalli and Dufwenberg (2007) propose a model of simple guilt, where players are assumed to be averse to letting down other players. More specifically, player $i$ 'lets down' player $j$ when his strategy provides player $j$ with a final payoff below the payoff expected by player $j$. Consider a game $\Gamma$ which satisfies condition H and a history $\tilde{h} \in \tilde{\mathcal{H}}$.

In history $\tilde{h}$ player i has two actions, $a'_{i,\tilde{h}}$ and $a''_{i,\tilde{h}}$. Define guilt from both actions as

$$B_i(a'_{i,\tilde{h}}, \mu_i^2(\tilde{h})) = \left[ \mathbf{E}_i(\mathbf{E}_j(\pi_j)) - \pi_j(a'_{i,\tilde{h}}) \right] 1 \left[ \pi_j(a'_{i,\tilde{h}}) \leq \pi_j(a''_{i,\tilde{h}}) \right] \tag{12}$$

$$B_i(a''_{i,\tilde{h}}, \mu_i^2(\tilde{h})) = \left[ \mathbf{E}_i(\mathbf{E}_j(\pi_j)) - \pi_j(a''_{i,\tilde{h}}) \right] 1 \left[ \pi_j(a'_{i,\tilde{h}}) > \pi_j(a''_{i,\tilde{h}}) \right] \tag{13}$$

where $1[A]$ denotes an indicator function taking a value of 1 when event $A$ occurs and 0 otherwise, $\mathbf{E}_j(\pi_j)$ denotes player $j$'s expectation of the own final payoff, and $\mathbf{E}_i(\mathbf{E}_j(\pi_j))$ denotes player $i$'s expectation of $\mathbf{E}_j(\pi_j)$. More formally

$$\begin{aligned} \mathbf{E}_i(\mathbf{E}_j(\pi_j)) &= \mathbf{E}_i(\mu_j^1(a'_{i,\tilde{h}}))\pi_j(a'_{i,\tilde{h}}) + (1 - \mathbf{E}_i(\mu_j^1(a'_{i,\tilde{h}})))\pi_j(a''_{i,\tilde{h}}) \\ &= \mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) \left[ \pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}}) \right] + \pi_j(a''_{j,\tilde{h}}) \end{aligned} \tag{14}$$

where $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) = \mathbf{E}_i(\mu_j^1(a'_{i,\tilde{h}}|\tilde{h}))$ and such that $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) \in [0,1]$. Assume without loss of generality that $\pi_j(a'_{i,\tilde{h}}) < \pi_j(a''_{i,\tilde{h}})$. From (13) it follows that player $i$ cannot feel guilt when choosing $a''_{i,\tilde{h}}$ given the later provides player $j$ with the highest of the two possible payoffs, i.e. $B_i(a''_{i,\tilde{h}}, \cdot) = 0$ for all $i$. On the other hand, player $i$ feels guilt when choosing $a'_{i,\tilde{h}}$ as this provides player $j$ with his lowest possible payoff. Hence,

$$\begin{aligned} \Delta B_{i,\tilde{h}} &= \left[ \mathbf{E}_i(\mathbf{E}_j(\pi_j)) - \pi_j(a'_{i,\tilde{h}}) \right] \\ &= \mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) \left[ \pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}}) \right] + \pi_j(a''_{j,\tilde{h}}) - \pi_j(a'_{i,\tilde{h}}) \end{aligned} \tag{15}$$

Inspection of (12), (13), and (15) reveals that condition I is satisfied in history $\tilde{h}$. Without knowledge of $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h})$, it follows that

$$\Delta B_{i,\tilde{h}} \in [0, \pi_j(a''_{i,\tilde{h}}) - \pi_j(a'_{i,\tilde{h}})] \tag{16}$$

where the lower bound $\underline{\Delta B_{i,\tilde{h}}} = 0$ is obtained by setting $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) = 1$, while the upper bound $\overline{\Delta B_{i,\tilde{h}}} = \pi_j(a''_{i,\tilde{h}}) - \pi_j(a'_{i,\tilde{h}})$ is obtained by setting $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h}) = 0$. It follows that all

games are of the type 4 presented in (6). This implies that the set $\mathcal{D}$ defined in Proposition 1 is empty and thus $\phi_\lambda^l = -\infty$. The conditional identification region of $\phi$ is then given by $[-\infty, \phi_\lambda^u]$, where

$$\phi_\lambda^u = \min_i \left[ \frac{Q\lambda - \Delta\pi_i}{\pi_j(a_{i,\tilde{h}}^{''}) - \pi_j(a_{i,\tilde{h}}^{'})} \right] \tag{17}$$

## 3.2 Example 2: reciprocity ($\phi \geq 0$)

Dufwenberg and Kirchsteiger (2004) propose a model of reciprocity where the psychological payoff of player $i$ in the last non-terminal history $\tilde{h}$, $B_i(a_{i,\tilde{h}}, \boldsymbol{\mu}_i^k)$, is given by the product $PK(\tilde{h}) \times K(a_{i,\tilde{h}})$. The first term $PK(\tilde{h})$ involves player $i$'s perception of the kindness of player $j$ towards him in history $\tilde{h}$. Let $\mathbf{E}_j\left(\pi_i|\tilde{h}\right)$ denote player $j$'s expectation of $i$'s final payoff in history $\tilde{h}$, and $\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right)$ denote player $i$'s expectation of $\mathbf{E}_j\left(\pi_i|\tilde{h}\right)$. That is

$$
\begin{aligned}
\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) &= \mathbf{E}_i(\mu_j^1(a_{i,\tilde{h}}^{'}|\tilde{h}))\pi_i(a_{i,\tilde{h}}^{'}) + (1 - \mathbf{E}_i(\mu_j^1(a_{i,\tilde{h}}^{'}|\tilde{h})))\pi_i(a_{i,\tilde{h}}^{''}) \\
&= \mu_i^2(a_{i,\tilde{h}}^{'}|\tilde{h}) \left[ \pi_i(a_{i,\tilde{h}}^{'}) - \pi_i(a_{i,\tilde{h}}^{''}) \right] + \pi_i(a_{i,\tilde{h}}^{''})
\end{aligned}
\tag{18}
$$

where $\mu_i^2(a_{i,\tilde{h}}^{'}|\tilde{h}) = \mathbf{E}_i(\mu_j^1(a_{i,\tilde{h}}^{'}|\tilde{h}))$. Moreover, define the 'equitable' payoff[7]

$$\pi_i^{e_j}(\mu_i^2(\tilde{h})) = \frac{1}{2} \left[ \max_{a_j \in A_j}\{\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i\right)\right)\} + \min_{a_j \in A_j}\{\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i\right)\right)\} \right].$$

The equitable payoff is used by player $i$ as a reference point to measure the kindness of player $j$ towards him. In particular, player $i$'s perceived kindness of player $j$ is given by the following difference

$$PK(\tilde{h}) = \mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) - \pi_i^{e_j}(\mu_i^2(\tilde{h}))$$

Expected payoffs $\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right)$ higher (lower) than the equitable payoff are thus perceived as kind (unkind). The second term entering the psychological payoff function

---

[7]For notational simplicity assume that all of player $j$'s strategies are *efficient* as defined by Dufwenberg and Kirchsteiger (2004) on p. 276.

involves the kindness of player $i$ towards player $j$ when choosing $a'_{i,\tilde{h}}$

$$
\begin{aligned}
K(a'_{i,\tilde{h}}) &= \pi_j(a'_{i,\tilde{h}}) - \frac{1}{2}\left[\pi_j(a'_{i,\tilde{h}}) + \pi_j(a''_{i,\tilde{h}})\right] \\
&= \frac{1}{2}\left[\pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}})\right]
\end{aligned}
$$

A similar expression follows for $K(a''_{i,\tilde{h}})$, the kindness when choosing $a''_{i,\tilde{h}}$. Multiplying $PK(\tilde{h})$ with $K(a'_{i,\tilde{h}})$ and $K(a''_{i,\tilde{h}})$ and rearranging gives

$$
B_i(a'_{i,\tilde{h}}, \mu_i^2(\tilde{h})) = \frac{1}{2}\left[\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) - \pi_i^{e_j}(\mu_i^2(\tilde{h}))\right]\left[\pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}})\right] \quad (19)
$$

$$
B_i(a''_{i,\tilde{h}}, \mu_i^2(\tilde{h})) = \frac{1}{2}\left[\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) - \pi_i^{e_j}(\mu_i^2(\tilde{h}))\right]\left[\pi_j(a''_{i,\tilde{h}}) - \pi_j(a'_{i,\tilde{h}})\right] \quad (20)
$$

Differencing (19) and (20) yields

$$
\Delta B_{i,\tilde{h}} = \left[\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) - \pi_i^{e_j}(\mu_i^2(\tilde{h}))\right]\left[\pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}})\right]
$$

Inspection of (19), (20), and (21) reveals that condition I is satisfied in history $\tilde{h}$. The values of $\underline{\Delta B_{i,\tilde{h}}}$ and $\overline{\Delta B_{i,\tilde{h}}}$ will depend on the signs of two terms,

$$
\left[\mathbf{E}_i\left(\mathbf{E}_j\left(\pi_i|\tilde{h}\right)\right) - \pi_i^{e_j}(\mu_i^2(\tilde{h}))\right], \text{ and } \left[\pi_j(a'_{i,\tilde{h}}) - \pi_j(a''_{i,\tilde{h}})\right],
$$

and thus will potentially vary across games. For games where the terms have the same sign, $\Delta B_{i,\tilde{h}}$ is a monotonically increasing function of $\mu_i^2(a'_{i,\tilde{h}}|\tilde{h})$. Without knowledge of $\mu_i^2(\tilde{h})$, it follows that

$$
\Delta B_{i,\tilde{h}} \in [\underline{\Delta B_{i,\tilde{h}}}, \overline{\Delta B_{i,\tilde{h}}}] \quad (21)
$$

where

$$
\underline{\Delta B_{i,\tilde{h}}} = \min_{\mu_i^2(\tilde{h})}\{\Delta B_{i,\tilde{h}}\} \quad (22)
$$

$$
\overline{\Delta B_{i,\tilde{h}}} = \max_{\mu_i^2(\tilde{h})}\{\Delta B_{i,\tilde{h}}\} \quad (23)
$$

Values in (22) and (23) can be used to estimate the endpoints using (9) and (10).

## 3.3 Estimation and inference on the bounds of the identification region

We propose a two step procedure to estimate the endpoints of the conditional identification region. We estimate $\lambda$ in a first step using the subset of games which satisfy condition

S by assuming that $\Delta\epsilon_{i,\tilde{h}}$ follows a logistic distribution. This is a standard logit model which can be estimated by Maximum Likelihood. Let $\widehat{\lambda}$ denote the estimated value of $\lambda$ obtained in this way. The second step consists of estimating the endpoints of the identification region conditional on the Maximum Likelihood estimate $\widehat{\lambda}$ and the estimated values of $Q_i$. We obtain estimates of $Q_i$ by inverting estimated choice probabilities derived from a reduced form model (equation (30) in the following section). Naïve estimators of the endpoints of the conditional identification region are given by the following sample counterparts to Proposition 1

$$\widehat{\phi}^l_{\widehat{\lambda}} = \max_{i \in \underline{I}} \widehat{\underline{\phi}}_i \tag{24}$$

$$\widehat{\phi}^u_{\widehat{\lambda}} = \min_{i \in \overline{I}} \widehat{\overline{\phi}}_i \tag{25}$$

where $\widehat{\underline{\phi}}_i$ and $\widehat{\overline{\phi}}_i$ are the estimated values of $\underline{\phi}_i$ and $\overline{\phi}_i$ defined in Proposition 1 (with the unknown value of $\lambda$ replaced with $\widehat{\lambda}$) and where $\underline{I}$ and $\overline{I}$ denote the set of games which can be used to estimate the lower and upper endpoints respectively. It is well known that the estimators (24) and (25) are possibly biased in finite samples. This reflects the fact that the expectation of the maximum (minimum) of random variables is generally higher (lower) than the maximum (minimum) of the expectations. We can thus expect $\widehat{\phi}^l_{\widehat{\lambda}}$ to have an upward finite sample bias while we can expect that $\widehat{\phi}^l_{\widehat{\lambda}}$ has a downward finite sample bias. This implies that naive estimators based on (24) and (25) will on average tend to produce overly narrow conditional identification regions.

Chernozhukov, Lee and Rosen (2009) (hereafter CLR) propose a median-unbiased estimator of the endpoints of the identification region and propose a method to construct confidence intervals which can take into account the two step nature of our approach. Here, we implement their proposed approach for parametric models (see their appendix C.1). In particular, we define

$$\widehat{\phi}^l_{\widehat{\lambda},\theta} = \max_{i \in \widehat{\underline{I}}} \left\{ \widehat{\underline{\phi}}_i - \widehat{\underline{G}}(\theta)s(i) \right\} \tag{26}$$

$$\widehat{\phi}^u_{\widehat{\lambda},\theta} = \min_{i \in \widehat{\overline{I}}} \left\{ \widehat{\overline{\phi}}_i + \widehat{\overline{G}}(\theta)s(i) \right\} \tag{27}$$

where $s(i)$ denotes the estimated standard error of either $\widehat{\underline{\phi}}_i$ or $\widehat{\overline{\phi}}_i$, $\widehat{\underline{G}}(\theta)$ denotes the

estimated $\theta-$quantile of $\max_{i\in\widehat{\underline{I}}}\left\{\left(\widehat{\underline{\phi}}_i-\underline{\phi}_i\right)/s(i)\right\}$, $\widehat{\overline{G}}(\theta)$ denotes the estimated $\theta-$quantile of $\min_{i\in\widehat{\overline{I}}}\left\{\left(\widehat{\overline{\phi}}_i-\underline{\phi}_i\right)/s(i)\right\}$, $\widehat{\underline{I}}$ and $\widehat{\overline{I}}$ denote estimated sets of games used to make inferences on the endpoints. Note that $\widehat{\underline{G}}(\theta)s(i)$ represents a bias correction term which intuitively enters negatively in (26) to correct for the upward bias of the estimator in (24). In a similar way, $\widehat{\overline{G}}(\theta)s(i)$ represents a bias correction term which enters positively in (27) to correct for the downward bias of the estimator in (25). Both $\widehat{\underline{G}}(\theta)s(i)$ and $\widehat{\overline{G}}(\theta)s(i)$ account for the sampling variability of $\widehat{\lambda}$ and $\hat{Q}_i$. Details concerning computation of $\widehat{\underline{G}}(\theta)s(i)$ and $\widehat{\overline{G}}(\theta)s(i)$ can be found in CLR.

Under conditions discussed in CLR it can be shown that

$$\lim_{N\to\infty}\Pr\left(\underline{\phi}\left(\lambda\right)\leq\widehat{\phi}^l_{\hat{\lambda},\theta}\right)=\theta \tag{28}$$

$$\lim_{N\to\infty}\Pr\left(\overline{\phi}\left(\lambda\right)\leq\widehat{\phi}^u_{\hat{\lambda},\theta}\right)=\theta \tag{29}$$

It follows that setting $\theta=0.5$ yields median-unbiased lower and upper endpoint estimators. These estimators are median-unbiased in the sense that the asymptotic probability that the estimated values lie above their true value is at least a half. Moreover, one sided $p\%$ confidence intervals can be obtained by computing $\widehat{\phi}^l_{\hat{\lambda},p}$ and/or $\widehat{\phi}^u_{\hat{\lambda},1-p}$ for the relevant endpoints. Finally, results in CLR imply that a valid $p\%$ confidence interval for $[\phi^l_\lambda,\phi^u_\lambda]$ can be obtained by computing $[\widehat{\phi}^l_{\hat{\lambda},p/2},\widehat{\phi}^u_{\hat{\lambda},1-p/2}]$.

# 4 Empirical illustration

## 4.1 Experimental design and data

Our experiment is based on the sequential game in Figure 1. The experiment was run in January and February 2010 via the LISS-panel, an Internet survey panel managed by CentERdata at Tilburg University. In total 2000 members of the panel were invited to participate in the experiment involving 500 payoffwise different games as shown in Figure 1. Only the associated monetary payoffs of the players differed across the games. The payoffs of the games randomly chosen from a set of similar games used in Bellemare,

Sebald, and Strobel (2010). Approximately 1/3 of the 500 payoffwise unique games were recoded to ensure that condition S holds, that is such that $\pi^j(l) = \pi^j(r)$ (dark circles in Figure 1).

Each panel member was initially randomly assigned a role and a payoffwise unique game in the following way. First, 1500 panel members were assigned the role of player $i$ while 500 panel members were assigned the role of player $j$. This role assignment allowed us to gather more decisions of $i$-players whose behavior is the primary focus of the paper. Subsequently, we randomly assigned each of the 500 payoff different games to three $i$ players and to one $j$-player. In other words, each of the 500 games could potentially be played by three $i$-players and one $j$-player.

Given the infrastructure of the LISS-panel, the game was played across two consecutive survey months. In the first month, only panel members assigned to the role of player $i$ were contacted and offered the possibility to participate in the experiment. Before revealing their role and specific game, they were provided general instructions, informed that 50 payoff-wise unique games would randomly be chosen ex-post and paid out two months later. Furthermore, they were given the possibility to withdraw from the experiment. After the revelation of their role and game, they were told that they would be making their decisions before $j$-players and that decisions would be matched ex-post. 1139 of the 1500 invited panel members accepted the invitation and completed the experiment in the role of player $i$.[8] Panel members who completed the experiment were first presented their unique game and then asked to send a message to player $j$. We allowed participants to send messages in order to increase their awareness concerning the other person they were grouped with. They could choose between two different messages and not sending a message:

| ☐ | *If you let me decide between l and r, I will choose l* |
| ☐ | *If you let me decide between l and r, I will choose r* |
| ☐ | I do not want to send a message |

---

[8]7 more invited panel members logged on but did not complete the experiment.

Each player $i$ then made his/her decision using the strategy method: $i$-players chose between $l$ and $r$ at history $h^1$ before knowing the decision of player $j$ at history $h^0$.

Panel members assigned to the role of player $j$ made their decisions during the second survey months. All $j$ players were first provided instructions and were informed that 50 payoff-wise unique games would randomly be chosen ex-post and paid out at the completion of the experiment. Again, before revealing their roles and games, they were given the possibility to withdraw from the experiment. 328 of the 500 invited panel members accepted the invitation and completed the experiment in the role of player $j$.[9] For the unique games for which we had more than one complete set of $i$-players decisions, we randomly chose one of them to be used in the interaction with player $j$. Invited panel members who accepted to participate in the experiment were then presented their unique game, were given the message of their matched $i$-player, and were asked to chose between $L$ and $R$ at history $h^0$ in the game.

After the second survey month we randomly chose 50 payoff-wise unique games (i.e. 15% of the 328 games that had been completed by one $i$ and one $j$ player) and paid the participants that had played these games according to the decisions that they had taken in the game. Average values of $\pi_i(R)$ and $\pi_j(R)$ were 28.386€ and 21.150€ respectively. Moreover, average values of $\pi_i(l)$ and $\pi_j(l)$ were 17.184€ and 25.899€ while corresponding averages of $\pi_i(l)$ and $\pi_j(l)$ were 18.746€ and 25.933€. Figure 5 illustrates the payoff variation of both players which follow from history $h^1$ in Figure 1. In particular, we plot $\Delta\pi_i = \pi_i(r) - \pi_i(l)$ and $\Delta\pi_j = \pi_j(r) - \pi_j(l)$ for all 500 randomly chosen games. Games for which condition S holds (i.e. $\Delta\pi_j = 0$) are denoted $Invariant$ and are marked by full circles. All other games are denoted $Variant$ and marked by empty circles. We can see that the payoff differences for player $j$ lie between -50€ and 50€ while payoff differences for player $i$ vary between -35€ and 35€.

Our data reveals that 70.45% of $j$ players (first movers) determined the final allocation by choosing the outside option. We perform a preliminary analysis of the decisions of $i$ players by estimating a logit model relating their decisions ($l$ or $r$ at history $h^1$) to the

---

[9] 7 more invited panel members logged on but did not complete the experiment.

difference in payoffs of both players as well as to their respective outside options. In particular, we estimate the following equation

$$\Pr(c = r|\Delta\pi_j, \Delta\pi_i, \pi_j(R), \pi_i(R)) =$$

$$F([\Delta\pi_i + \alpha_1\Delta\pi_j + \alpha_2\pi_j(R) + \alpha_3\pi_i(R)]/\tilde{\lambda}). \tag{30}$$

where (30) can be interpreted as a reduced form model of equation (2). We find that the probability that $i$ players chooses $r$ increases significantly with $\Delta\pi_j$ ($\hat{\alpha}_1 = 0.160$, se. $= 0.043$), suggesting that $i$ players take into account the well being of $j$ players. Not surprisingly, the size of $\hat{\alpha}_1$ is substantially lower than 1, an indication that $i$ players value their own well-being more than that of others. Interestingly, we do not find that any of the outside options have a significant impact on the decisions of $i$ players ($\hat{\alpha}_2 = 0.103$, $p$-value $= 0.221$; $\hat{\alpha}_3 = $ -0.006, $p$-value $= 0.928$). Finally, we estimated an extended specification where we allowed the noise parameter $\tilde{\lambda}$ to depend on $\Delta\pi_i$ and $\Delta\pi_j$ by specifying $\tilde{\lambda} = \exp(\gamma_0 + \gamma_1\Delta\pi_i + \gamma_2\Delta\pi_j)$. We found no significant increase in the log-likelihood function value ($p$-value $= 0.9531$), suggesting that the noise level does not vary with the level of payoff differences of each player in the game. This provides some indication that condition N is likely to hold in the data.

## 4.2  Results for guilt aversion

Consider first the model of guilt aversion discussed in section 3.1 in the context of Figure 1, i.e. the strategic environment underlying our experiment. Furthermore, denote by $l$ the action of player $i$ which implies the higher payoff for player $j$, i.e. $\pi_j(r) < \pi_j(l)$. Given this the conditional identification region of $\phi$ is given by $[-\infty, \phi_\lambda^u]$, where

$$\phi_\lambda^u = \min_i \left[ \frac{Q_i\lambda - \Delta\pi_i}{\pi_j(l) - \pi_j(r)} \right] \tag{31}$$

We first assess what can be learned about the model parameters without exploiting the invariance condition. The grey area in Figure 2 presents the estimated identification region for $(\phi, \lambda)$ derived by computing (31) replacing $Q_i$ with $\hat{Q}_i$ for different values of $\lambda$. The diagonal line presents the locus of values of $\phi_\lambda^u$ for a selected range of values

of $\lambda$. We see that $\phi_\lambda^u$ is below zero for values of $\lambda$ between 0 and (approximately) 21, suggesting that players are guilt averse over this range of $\lambda$ values. However, $\phi_\lambda^u$ equals zero when $\lambda$ is greater than 21. It follows that the data is largely uninformative about the relevance of guilt aversion in our experiment when we do not exploit the invariance condition. This is one illustration of the analysis in Manski (2010) where he argues that choice data alone are in general insufficient to make meaningful inferences on preferences without information about the beliefs of players.

We next exploited the invariance condition to estimate (31) by replacing $\lambda$ with a consistent estimate obtained in a first step using games which satisfy condition S. Table 1 presents the results. Column $\lambda$ contains the estimated value of the scale parameter while column $\left(-\infty, \widehat{\phi}_{\widehat{\lambda}}^u\right]$ presents the estimated identification region using the naive endpoint estimator based on (17). As discussed in section 3.3, the naive estimator is potentially biased downwards in finite samples. Columns $\widehat{\phi}_{\widehat{\lambda},0.5}^u$ and $\widehat{\phi}_{\widehat{\lambda},0.95}^u$ present the median-unbiased estimator and the corresponding one-sided 95% confidence band based on CLR.

The estimated value of $\lambda$ obtained using all games which satisfy condition S is 14.140 and is significant at the 1% level.[10] This estimate implies that $\phi_{\widehat{\lambda}}^u$ is estimated to be -0.881, suggesting that players are on average willing to pay at least 0.88€ to avoid letting down player $j$ by 1€. This value can alternatively be derived from Figure 2 which plots $\widehat{\lambda}$ and the corresponding estimated values of $\phi_{\widehat{\lambda}}^u$. Inspection of the Figure illustrates the identification power of conditions I and S – the identification region is reduced to a single (vertical) line. Column $\widehat{\phi}_{\widehat{\lambda},0.5}^u$ reveals that the downward bias of these estimated upper endpoints is substantial. In particular, the estimated upper endpoint for the entire sample

---

[10]This suggests that a significant proportion of $i$ players chose the option providing them with the lowest payoff, given the payoff invariance for player $j$. One interpretation of this result is that $\Delta\epsilon_{i,\tilde{h}}$ captures noise and sub-optimal decision making. Another interpretation is that part of $\Delta\epsilon_{i,\tilde{h}}$ captures unobserved preferences such as inequity aversion. Then, some players may be selecting the lowest payoff for themselves in order to reduce the payoff difference with player $j$. This would be consistent with results presented in Bellemare, Kröger, and van Soest (2008) who analyze responder behavior in the ultimatum game in the Dutch population. They found that a substantial proportion of responders were willing reject overly generous offers which provided them higher payoffs than proposers.

increases from -0.881 to -0.475 when controlling for the finite sample bias. The last column of the table presents the estimated one-sided 95% confidence interval for $\phi_\lambda^u$. Values less than zero reveal significant guilt aversion. The estimated 95% confidence interval for $\phi_\lambda^u$ is -0.077, suggesting significant guilt aversion in the broad population.

We then repeated the analysis for different sub-groups of the population. In particular, we performed a separate analysis for men and women, for three education levels (low, intermediate, and high levels), and for two age groups (below or above median sample age). Finer partitions potentially including other socio-economic variables or their interactions are in principle possible. However, our chosen partitions ensure that we have sample sizes which allow us to make meaningful comparisons. We find that the estimated values of $\lambda$ are positive and significant at the 1% level for all sub-populations considered. The estimated values of $\phi_\lambda^u$ vary substantially across the sub-populations. For example, players with low education levels have the highest estimated upper endpoint (-0.337) while players with high levels of education have the lowest estimated upper endpoint (-1.306). The bias-corrected estimated upper endpoints for the other partitions are also substantially higher then the corresponding estimates based on the naive estimator, suggesting important finite sample bias for the naive endpoint estimator. Overall, the median bias-corrected upper endpoints vary from -0.871 (men) to 0.029 (low education). Finally, the estimated one-sided 95% confidence intervals for $\phi_\lambda^u$ suggest that guilt aversion is significant for men, high educated players, and players above 47 years of age.

## 4.3   Results for reciprocity

We now consider the possibility that players have reciprocal preferences as outlined in section 3.2. Given the strategic environment displayed in Figure 1 the equitable payoff $\pi_i^{e_j}$ and the perceived kindness $PK(h^1)$ of player $i$ in history $h^1$ respectively reduce to

$$\pi_i^{e_j} = \frac{1}{2} \left[ \mathbf{E}_i \left( \mathbf{E}_j \left( \pi_i | h^1 \right) \right) + \pi_i(R) \right],$$

with

$$\mathbf{E}_i \left( \mathbf{E}_j \left( \pi_i | h^1 \right) \right) = \mathbf{E}_i(\mu_j^1(r|h^1))\pi_i(r) + (1 - \mathbf{E}_i(\mu_j^1(r|h^1)))\pi_i(l)$$

$$= \mu_i^2(r|h^1) \left[ \pi_i(r) - \pi_i(l) \right] + \pi_i(l)$$

and

$$PK(h^1) = \mathbf{E}_i \left( \mathbf{E}_j \left( \pi_i | h^1 \right) \right) - \frac{1}{2} \left[ \mathbf{E}_i \left( \mathbf{E}_j \left( \pi_i | h^1 \right) \right) + \pi_i(R) \right]$$

$$= \frac{1}{2} \left[ \mu_i^2(r|h^1) \left[ \pi_i(r) - \pi_i(l) \right] + \pi_i(l) - \pi_i(R) \right]$$

Table 2 presents the results for the same sub-populations used in our analysis of guilt aversion. All results concerning the estimation of $\lambda$ are identical to the one presented for guilt aversion. Column $[\widehat{\phi}_{\hat{\lambda}}^l, \widehat{\phi}_{\hat{\lambda}}^u]$ presents the identification region estimated using the naive endpoint estimator for both endpoints. Columns $\widehat{\phi}_{\hat{\lambda},0.5}^l$ and $\widehat{\phi}_{\hat{\lambda},0.025}^l$ present respectively the median-unbiased estimated lower endpoint and the corresponding one-sided 97.5% confidence band using the approach proposed by CLR. Columns $\widehat{\phi}_{\hat{\lambda},0.5}^u$ and $\widehat{\phi}_{\hat{\lambda},0.975}^u$ present the corresponding estimates for the upper endpoint of the identification region. The interval $\widehat{\phi}_{\hat{\lambda},0.025}^l, \widehat{\phi}_{\hat{\lambda},0.975}^u$ forms a 95% confidence interval for the identification region $\phi_{\lambda}^l, \phi_{\lambda}^u$.

We find that the naive estimator produces estimated endpoints which cross: the estimated values of $\phi_{\hat{\lambda}}^l$ exceed the estimated values of $\phi_{\hat{\lambda}}^u$ for all sub-populations considered.[11] Moreover, the estimated upper endpoints are censored at zero for all sub-populations. Both these results can be explained by the fact that naive estimators of the lower (upper) endpoints are potentially biased upwards (downwards) in finite samples. We find that the median-unbiased estimator of CLR resolves most of the crossings observed when using the naive estimators. A notable exception concerns the sub-population of players with intermediate levels of education. There, the median-unbiased estimated lower endpoint remains slightly above the median-unbiased estimated upper endpoint. We find that all 95% confidence intervals $[\widehat{\phi}_{\hat{\lambda},0.025}^l, \widehat{\phi}_{\hat{\lambda},0.975}^u]$ are narrow and are either close to zero or overlap with zero. In line with the aforementioned fact that outside options did not have a

---

[11]Crossing of endpoints estimated using "naive" estimators of the form discussed in this paper are not uncommon. Chesher (2009) provides further examples.

significant impact on the decisions of $i$ players, these results suggests that reciprocity is a weak predictor of decisions of $i$ players in our experiment.

# 5 Conclusion

We proposed an approach to learn about the empirical relevance of belief-dependent preferences in sequential two-player games without exploiting information about the beliefs of players. Our approach exploits the natural bounds of the psychological payoffs of players to make set inferences on the relevance of the underlying belief-dependent preferences. Existing research has established that the identification regions of the model parameters are typically large an uninformative without information on beliefs. However, we showed that the identification regions can be substantially reduced by exploiting a simple invariance property which is embedded in several prominent belief-dependent preferences.

Our approach produced informative bounds for the relevance of belief-dependent preferences in our experiment. In particular, our analysis of guilt aversion suggests that the population willingness to pay to avoid letting down the other player by 1€ is significantly different from zero and at least greater or equal to 0.08€. We also found that several groups of the population are willing to pay more at a minimum. In particular, high educated individuals are willing to pay at least 0.14€ while men are willing to pay at least 0.16€ to avoid letting down the other player by 1€. We were also able to obtain tight and very informative bounds around the relevance of reciprocity in our experiment. Our results suggest that reciprocity weakly predicts the final decisions made in our experiment. This result holds for all groups of the population we considered.

These results can be interpreted as providing approximate bounds around the *average* sensitivity parameter for each of the sub-groups of the population considered. Researchers may additionally want to conduct an individual-specific analysis to learn about the entire distribution of the sensitivity parameter within each sub-group of the population. Our approach can in principle be extended to make individual-specific inferences by exploiting data from subjects making multiple decisions in games satisfying condition S and games

where payoffs of player $j$ vary with the action taken by player $i$. Future work should determine the properties of the proposed approach in relationship to the number of decisions available for each subject in order to bound individual sensitivity parameters. Future work should also try to extend the approach to settings with more than two decisions as well as to settings where researchers are interested in combining data from different games.

The approach proposed in this paper ultimately allows researchers to assess the added value of exploiting data on stated beliefs to learn about the relevance of belief-dependent preferences in games. Our analysis of reciprocity provides an example where little can be gained by further exploiting stated belief-data: the estimated identification regions are narrow and precisely estimated. Our results also suggests that this result is unlikely to hold in general. Estimated identification regions in the case of guilt aversion remain large despite revealing significant guilt aversion in various sub-groups of the population. Researchers requiring more precise information about the exact level of guilt aversion (or other preferences in the class) must then exploit data on higher-order beliefs to point identify the sensitivity parameters. This will require more work to carefully address the possibility that stated beliefs are measured with error and/or correlated with preferences entering the model.
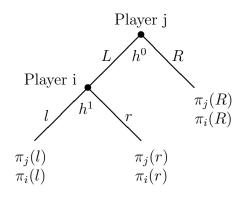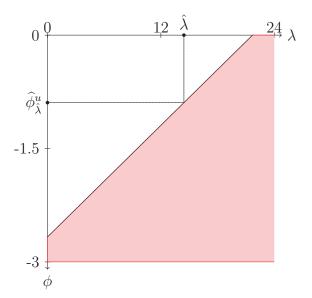
Figure 1: The Game



Figure 2: Estimated identification region for $(\phi, \lambda)$ in the case of simple guilt. $\hat{\lambda}$ denotes the value of $\lambda$ estimated using all games which satisfy condition S. $\widehat{\phi}^u_{\hat{\lambda}}$ denotes the estimated upper bound of the identification region of $\phi$.
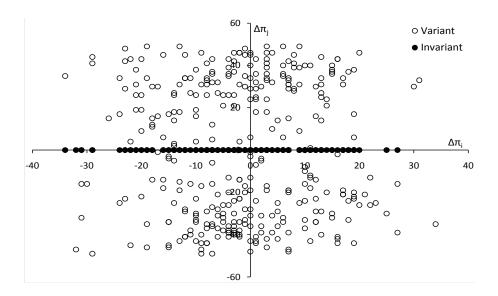
Figure 3: Payoff variation across games for both roles, with $\Delta\pi_i = \pi_i(l) - \pi_i(r)$ on the horizontal axis and $\Delta\pi_j = \pi_j(l) - \pi_j(r)$ on the vertical axis. "Invariant" denotes games where $\Delta\pi_j = 0$ to satisfy condition S. Other games (empty circles) are denoted as "Variant".

| Guilt aversion | $N_1$ | $N_2$ | $\widehat{\lambda}$ | $(-\infty, \widehat{\phi}^u_\lambda]$ | $\widehat{\phi}^u_{\lambda,0.5}$ | $\widehat{\phi}^u_{\lambda,0.95}$ |
|---|---|---|---|---|---|---|
| All sample | 349 | 648 | 14.440 | $(-\infty, -0.881]$ | -0.475 | -0.077 |
| Women | 184 | 344 | 18.198 | $(-\infty, -0.465]$ | -0.064 | 0.067 |
| Men | 165 | 304 | 11.013 | $(-\infty, -1.287]$ | -0.871 | -0.159 |
| Low education | 109 | 192 | 20.467 | $(-\infty, -0.337]$ | 0.029 | 0.228 |
| Intermediate education | 126 | 242 | 12.347 | $(-\infty, -0.755]$ | -0.179 | 0.032 |
| High education | 114 | 214 | 13.223 | $(-\infty, -1.305]$ | -0.726 | -0.144 |
| Age $\leq 47$ | 182 | 324 | 17.634 | $(-\infty, -0.550]$ | -0.154 | -0.002 |
| Age $> 47$ | 167 | 324 | 12.044 | $(-\infty, -0.987]$ | -0.536 | -0.047 |

Table 1: Results of the two step procedure when assuming that players are guilt averse. The table presents the sample sizes used in both estimation steps. $N_1$ denotes the sample size to estimate $\lambda$ in the first step. $N_2$ denotes the sample size to estimate $[-\infty, \phi^u_\lambda]$. The column $(-\infty, \phi^u_\lambda]$ presents the estimated identification region based on the naive estimator. Columns $\widehat{\phi}^u_{\lambda,0.5}$ and $\widehat{\phi}^u_{\lambda,0.95}$ present respectively the median-unbiased estimated upper bound and the one-sided 95% confidence band based on Chernozhukov, Lee and Rosen (2009). All estimated value of $\lambda$ reported in the table are significant at the 1% level.

28

| *Reciprocity* | $N_1$ | $N_2$ | $\widehat{\lambda}$ | $\widehat{\phi}^l_{\lambda,0.025}$ | $\widehat{\phi}^l_{\lambda,0.5}$ | $[\widehat{\phi}^l_\lambda, \widehat{\phi}^u_\lambda]$ | $\widehat{\phi}^u_{\lambda,0.5}$ | $\widehat{\phi}^l_{\lambda,0.975}$ |
|---|---|---|---|---|---|---|---|---|
| All sample | 349 | 648 | 14.440 | 0.006 | 0.016 | [0.067, 0.000] | 0.020 | 0.031 |
| Women | 184 | 344 | 18.198 | -0.003 | 0.012 | [0.031, 0.000] | 0.027 | 0.041 |
| Men | 165 | 304 | 11.013 | 0.005 | 0.042 | [0.102, 0.000] | 0.012 | 0.027 |
| Low education | 109 | 192 | 20.467 | -0.018 | -0.001 | [0.026, 0.000] | 0.031 | 0.048 |
| Intermediate education | 126 | 242 | 12.347 | 0.005 | 0.022 | [0.039, 0.000] | 0.014 | 0.029 |
| High education | 114 | 214 | 13.223 | 0.004 | 0.019 | [0.051, 0.000] | 0.027 | 0.047 |
| Age $\leq$ 47 | 182 | 324 | 17.634 | -0.001 | 0.017 | [0.045, 0.000] | 0.022 | 0.038 |
| Age $>$ 47 | 167 | 324 | 12.044 | 0.003 | 0.014 | [0.048, 0.000] | 0.016 | 0.033 |

Table 2: Results of the two step procedure when assuming that players have potentially reciprocal preferences. The table presents the sample sizes used in both estimation steps. $N_1$ denotes the sample size to estimate $\lambda$ in the first step. $N_2$ denotes the sample size to estimate $[\phi^l_\lambda, \phi^u_\lambda]$. The column $[\widehat{\phi}^l_\lambda, \widehat{\phi}^u_\lambda]$ presents the estimated identification region based on the naive estimator. Columns $\widehat{\phi}^l_{\lambda,0.5}$ and $\widehat{\phi}^u_{\lambda,0.025}$ present respectively the median-unbiased estimated lower bound and the corresponding one-sided 95% confidence bands based on Chernozhukov, Lee and Rosen (2009). Columns $\widehat{\phi}^u_{\lambda,0.5}$ and $\widehat{\phi}^l_{\lambda,0.975}$ present the corresponding estimates for the upper bound of the identification region. All estimated value of $\lambda$ reported in the table are significant at the 1% level.

# References

BATTIGALLI, P., AND M. DUFWENBERG (2007): "Guilt in Games," *American Economic Review Papers and Proceedings*, 97, 170–176.

——— (2009): "Dynamic Psychological Games," *Journal of Economic Theory*, 144, 1–35.

BELLEMARE, C., L. BISSONNETTE, AND S. KRÖGER (2010): "Bounding Preference Parameters under Different Assumptions about Beliefs: a Partial Identification Approach," *Experimental Economics*, 13, 334–345.

BELLEMARE, C., AND S. KRÖGER (2007): "On Representative Social Capital," *European Economic Review*, 51, 183–202.

BELLEMARE, C., S. KRÖGER, AND A. VAN SOEST (2008): "Measuring Inequity Aversion in a Heterogeneous Population using Experimental Decisions and Subjective Probabilities," *Econometrica*, 76, 815–839.

BELLEMARE, C., A. SEBALD, AND M. STROBEL (2011): "Measuring the Willingness to Pay to Avoid Guilt: Estimation using Equilibrium and Stated Belief Models," *forthcoming, Journal of Applied Econometrics*.

BLANCO, M., D. ENGELMANN, A. KOCH, AND H.-T. NORMANN (2011): "Preferences and Beliefs in a Sequential Social Dilemma: a Within-subjects Analysis," *Working paper, Manheim University*.

CAPPELAN, A., A. HOLE, E. SØRENSEN, AND B. TUNGODDEB (2007): "The Pluralism of Fairness Ideals: An Experimental Approach," *American Economic Review*, 97, 818–827.

CHARNESS, G., AND M. DUFWENBERG (2006): "Promises and Partnerships," *Econometrica*, 74, 1579–1601.

——— (2010): "Participation," *American Economic Review, forthcoming*.

CHERNOZHUKOV, V., S. LEE, AND A. ROSEN (2009): "Intersection Bounds: Estimation and Inference," *Working paper*.

CHESHER, A. (2009): "Single equation endogenous binary response models," *Cemmap working paper 23/09*.

DHAENE, G., AND J. BOUCKAERT (2010): "Sequential reciprocity in two-player, two-stage games: An experimental analysis," *Games and Economic Behavior*, 70, 289–303.

DUFWENBERG, M., AND G. KIRCHSTEIGER (2004): "A Theory of Sequential Reciprocity," *Games and Economic Behavior*, 47, 268–298.

FALK, A., E. FEHR, AND U. FISCHBACHER (2008): "Testing theories of fairness-Intentions matter," *Games and Economic Behavior*, 62, 287–303.

GEANAKOPLOS, J., D. PEARCE, AND E. STACCHETTI (1989): "Psychological Games and Sequential Rationality," *Games and Economic Behavior*, 1, 60–79.

KLEIN, R. W., AND R. H. SPADY (2002): "An Efficient Semiparametric Estimator for binary Response Models," *Econometrica*, 61, 387–421.

MANSKI, C. (2010): "Random Utility Models with Bounded Ambiguity," in *Structural Econometrics, Essays in Methodology and Applications*, ed. by D. Butta, pp. 272–284. Oxford University Press, New Delhi.

MANSKI, C. F., AND E. TAMER (2002): "Inference on Regressions with Interval Data on a Regressor or Outcome," *Econometrica*, 70, 519–546.

VANBERG, C. (2010): "A Short Note of the Rationality of the False Consensus Effect," *Mimeo, Heidelberg University*.