



Is the Veil of Ignorance Transparent ?

Gaël Giraud, Cécile Renouard

► To cite this version:

| Gaël Giraud, Cécile Renouard. Is the Veil of Ignorance Transparent ?. 2010. halshs-00593973

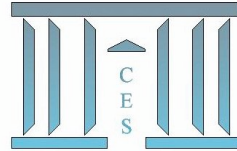
HAL Id: halshs-00593973

<https://shs.hal.science/halshs-00593973>

Submitted on 18 May 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Is the Veil of Ignorance Transparent ?

Gaël GIRAUD, Cécile RENOARD

2011.26



Is the Veil of Ignorance Transparent?*

GAËL GIRAUD and CÉCILE RENOUARD

April 19, 2011

ABSTRACT.— Theories of justice in the spirit of Rawls and Harsanyi argue that fair-minded people should aspire to make choices for society as if in the original position, that is, behind a veil of ignorance that prevents them from knowing their own social positions. In this paper, we provide a fairly simple framework showing that preferences in front of the veil of ignorance (i.e., in face of everyday risky situations) can be entirely deduced from ethical preferences behind the veil. Moreover, by contrast with Kariv & Zame (2008), in many cases of interest, the converse is not true: Ethical decisions cannot be deduced from economic ones. This not only rehabilitates distributive theories of justice but even proves that standard decision theory in economic environments cannot be exonerated from ethical questioning.

KEYWORDS. Moral preferences, business ethics, social preferences, distributional justice, theory of justice, social choice, original position, veil of ignorance, utilitarianism, maximin principle.

JEL Classification: D63.

* We thank, without implicating, Jacques Drèze and François Maniquet, as well as two anonymous referees and the participants of the “Cercle d’Epistémologie économique” (University Paris-1) for helpful comments. The usual *caveat* applies.

1 Introduction

Rawls (1971, 1974) and Harsanyi (1953, 1955, 1975) have constructed theories of social justice based on the choices that representatives should make for society in what Rawls names the “original position”, behind a veil of ignorance that prevents people from knowing their own future positions. Rawls (1971) views preferences in the original position as having a different nature from “ordinary” preferences for consumption, for risk or for the distribution of social goods to others. Rawls specifies that the parties in the original position are concerned only with citizens’ share of what he calls primary social goods, which include basic rights as well as economic and social advantages. Rawls also argues that the representatives in the original position would adopt the maximin rule as their principle for evaluating the choices before them, i.e., making the choice that produces the highest payoff for the least advantaged position. Being behind the veil of ignorance guarantees that the conception of justice to emerge will be agreed upon in a fair situation. “Fairness of the circumstances under which agreement is reached transfers into the fairness of the principles agreed to” (Rawls (1974)). Since these principles serve as principles of justice, the veil of ignorance therefore plays a crucial role in Rawls’ construction of “justice as fairness”.

In this paper, we shall consider the situation of a Representative who can face three types of decision-making problem: (1) Behind the veil of ignorance, her preferences will be called “ethical”;¹ (2) in a risky individual decision problem (in front of the veil), her preferences will be termed “risk preferences”; (3) finally, in a social choice problem (still in front of the veil, since the Representative is assumed to know her position), her preferences will be “social”.²

If Rawls and Harsanyi come to quite different conclusions about the form ethical preferences should take behind the veil of ignorance respectively the maximin and the “utilitarian” criteria, this is mainly due to their different view on the attitude of people towards uncertainty behind the veil of ignorance. Nevertheless, both Harsanyi and Rawls agree to view the original position as a purely hypothetical situation, a thought experiment where ethical preferences are theoretical constructs that should conform to some rationality requirements, paving the road towards various theories of justice.

By contrast, Kariv and Zame (2008) have recently introduced a framework encompassing both risk, social and ethical preferences, where they show that, under some assumptions, ethical preferences in the original position are entirely determined by risk and social preferences, i.e., by preferences that are not hypothetical at all. In other words, according to these authors, preferences behind the veil of ignorance can be deduced from preferences in front of the veil of ignorance. Since these authors view risk and social preferences as being essentially arbitrary, they conclude that “there is no conceptual reason to expect that moral preferences should be consistent with any particular notion of rationality or theory of justice”. Thus, at variance with both Rawls and Harsanyi, Kariv and Zame (2008) reach a conclusion similar

¹Following Ricœur’s (1992, p.170) distinction between ethics as the aim of an accomplished life (teleological perspective) and morality as the norms related to a deontological point of view, we prefer here the term “ethical” to “moral”.

²For a recent survey of the literature on social preferences, see Fehr and Schmidt (2006).

to that of Hayek (1976), according to whom social justice is a “mirage”.

In this paper, we challenge this viewpoint by reexamining the framework introduced by Kariv and Zame (2008). Starting with the same setting, but adopting different assumptions (which encompass more classic preferences than do the assumptions needed by Kariv and Zame (2008)), we provide an extremely simple proof of exactly the opposite result: Risk and social preferences can be entirely deduced from ethical ones. Moreover, we show by means of examples (see subsection 4.2 below) that many cases of interest (such as the leximin criterion or utilitarianism) do not fulfill Kariv and Zame (2008) assumptions but verify the axioms of this paper. In such examples, not only do risk and social preferences follow from ethical ones, but the converse is not true: Ethical preferences cannot be deduced from risk and social ones. Thus, we agree with Kariv and Zame (2008) that there is a link between preferences behind and in front of the veil of ignorance. In our view, however, the implication goes in the reverse direction: Theories of justice cannot be reduced to descriptive theories (how people actually behave *de facto*) but are indeed normative theories (how people ought to choose). As for risk and social preferences, they cannot be reduced to descriptive rules of thumb either: They belong to prescriptive theories (i.e., practical aids to choice) which follow from ethical decisions.

The next section provides our set-up. Section 3 shows how risk and social preferences can be deduced from ethical ones. Section 4 studies the reverse relationship. Examples and counter-examples are provided in section 5. Finally, the last section contains a discussion of the main conceptual issues at stake.

2 Choice environments

Following Kariv & Zame (2008), society consists of N actors, $i = 1, \dots, N$, of whom there is no loss of generality in assuming that the Representative is player 1.³ Three environments are considered. In the first, termed the **ETHICAL CHOICE** environment, the objects of choice are allocations of prospects for all members of the society, including the Representative, but in a setting where the Representative does not know her position in the society, nor the positions of others. In the second environment, called the **SOCIAL CHOICE** environment, the objects of choice are (deterministic) allocations of prospects for all the members of society, including the Representative. By contrast with the **ETHICAL** environment, the Representative in the social choice environment knows what her social position will be before taking a decision. In the third, which we term the **RISK** environment, objects of choice are random individual prospects for the Representative. As for prospects, they may designate a huge variety of items: utility levels, income, poverty indices, etc.⁴

Choice spaces are formalized as follows:

- The choice space \mathcal{R} in the individual risk environment consists of all lotteries, that is, collections

$$(p_j x_j)_{j=1, \dots, K}, \quad (1)$$

³In subsection 4.2, Example 2, below, an alternate interpretation of the indices $i = 1, \dots$ will be proposed.

⁴Notice that preferences need not be increasing with respect to prospects.

where $(p_j)_j$ is a probability vector⁵, and each $x_j \in \mathbb{R}$ is a prospect.⁶ The lottery (1) yields the Representative prospect x_j with probability p_j .

- The choice space \mathcal{S} in the social choice environment consists of all deterministic allocations x^j (not to be confused with x_j) in \mathbb{R}^N . This allocation yields the citizen $i \in \mathbb{N}$ the prospect x_i^j with certainty.

Let $\text{Perm}(N)$ be the group of permutations $\sigma : N \rightarrow N$. Given some vector $x \in \mathbb{R}^N$ and some permutation $\sigma \in \text{Perm}(N)$, the composition $x\sigma$ is again an element of \mathbb{R}^N assigning prospect $x_{\sigma(i)}$ to individual i .

- The choice space \mathcal{E} in the ethical choice environment consists of all lotteries

$$(p_\sigma(x\sigma))_{\sigma \in \text{Perm}(N)}$$

where (p_σ) is a probability distribution on the set $\text{Perm}(N)$ and $x \in \mathbb{R}^N$. This lottery yields citizen i prospect $x_{\sigma(i)}$ with probability p_σ . In particular, it provides the Representative with prospect $x_{\sigma(1)}$ with probability p_σ (for all $\sigma \in \text{Perm}(N)$).

In the risk environment, the Representative is simply a decision maker who must choose a random prospect for herself. In the social choice environment, the Representative is to choose a deterministic prospect for every individual in the society. In the ethical choice environment, she is to choose a deterministic distribution of prospects across society but with the random assignment of individuals to places in society. When interpreted in terms of Knightian riskiness, this ethical choice environment (with equal probabilities) coincides with Harsanyi's (1953,1955) formalization of ethical decisions (more on this in Example 2 *infra*).

One consequence of our approach is that the set of ethical choices is rich enough to include all social choices. On *this* aspect, we agree with Kariv & Zame (2008): If part of the social choice problem was to escape from any ethical questioning, our standpoint is that this would have to be a priori justified.

In the ethical, social choice and risk environments, the Representative's preference relations are written \succeq_e, \succeq_s and \succeq_r respectively. In order to be able to shift from one environment to the other, we need to consider a global set-up encompassing both \mathcal{E}, \mathcal{S} and \mathcal{R} . Let us therefore denote by \mathcal{L} the space of lotteries over allocations:

$$(p_j x^j)_j, \tag{2}$$

where each $x^j \in \mathbb{R}^N$ is an allocation. The lottery given by (2) yields x_1^j to the Representative with probability p_j . Sometimes, we write (2) in the form $(p_j x^j)$ when the index is clear. $(1x^j)$ stands for the degenerate lottery yielding allocation x^j for sure.

⁵That is, $p_j \geq 0$, for each j , and $\sum_j p_j = 1$.

⁶A prospect may be an income, a utility level or any quantitative characterization of an economic situation. For simplicity, they are assumed, here, to be real numbers but prospects might take value in a multi-dimensional space (or an abstract space of outcomes) without impairing our results.

The space \mathcal{L} encompasses the three environments mentioned supra. To see that $\mathcal{R} \subset \mathcal{L}$, it suffices to identify the individual lottery $(p_j x_j) \in \mathcal{R}$ with the lottery of collective allocations $(p_j(x_j, 0, \dots, 0)) \in \mathcal{L}$. That is, identify \mathcal{R} as the subset of \mathcal{S} consisting of lotteries that yield all individuals but the Representative the 0 prospect almost surely. There is no doubt that this identification is highly disputable. We borrow it from Kariv & Zame (2008). The discussion of alternative ways to tackle this issue would go beyond the scope of this paper and is postponed to further work. Similarly, \mathcal{S} is identified with the subset of \mathcal{L} consisting of degenerate lotteries ; \mathcal{E} is identified with the subset of \mathcal{L} consisting of lotteries of the form $(p_\sigma x^\sigma)_\sigma$ with the property that $x^\sigma = x\sigma$ for each $\sigma \in \text{Perm}(N)$.

3 Deducing risk preferences from ethical preferences

Preferences will be characterized by two postulates. The first gathers hardly controversial rationality requirements on global preferences \succeq over \mathcal{L} . Before stating them explicitly, let us recall the definition of compound lotteries. Suppose that L_1, \dots, L_K are K lotteries, and $(p_k)_{k=1, \dots, K}$ is a probability distribution. Then, $(p_k L_k)_k$ denotes a compound lottery in the following sense: One and only one lottery will be the prize, and the probability that it will be L_k is p_k .

A0 (i) Transitivity. The relation \succeq on \mathcal{L} is transitive.

(ii) Reduction of compound lotteries. Any compound lottery in \mathcal{L} is indifferent to a simple lottery, their probabilities being computed according to the ordinary calculus. In particular, if $(q_k)_k$ is a probability distribution and each $L_k = (p_k^i x^i)_i$ for $k = 1, \dots, K$, is a lottery, then there is no loss of generality in assuming that they all involve the same finite set, $(x_j)_j$, of allocations, and moreover

$$(q_k L_k)_k \sim (\tilde{p}_j x^j)_j$$

with $\tilde{p}_j := \sum_k q_k p_k^j$.

(iii) Continuity. Given any collection of allocations $(x_1, \dots, x_K) \in \mathcal{S}^K$, ordered (without loss of generality) so that $x^i \succeq_s x^{i+1}$ for every $1 \leq i \leq K-1$, then every x_i is indifferent in \mathcal{L} to some lottery involving only x_1 and x_K , i.e., there exists a probability $p_i \in [0, 1]$ such that

$$x^i \sim (p_i x^1, 0x^2, \dots, 0x^{K-1}, (1-p_i)x^K) =: X^i. \quad (3)$$

(iv) Substitutability. In any lottery $(p_k x^k)_k$ and for every i , X^i (as defined by (3)) can be substituted to x^i , that is: $(p_1 x^1, \dots, p_K x^K) \sim (p_1 x^1, \dots, p_i X^i, \dots, p_K x^K)$.

A0 (iii) is a continuity assumption on global preferences \succeq .⁷ Suppose, indeed, that $x^1 \succ_s x^2 \succ_s x^3$. It is plausible that the lottery $(p x^1, (1-p)x^3)$ is preferred

⁷Cf. Luce and Raiffa (1957), p. 27.

to x^2 as p approaches 1, and that the preference is inverted when p is close to 0. This assumption simply says that, as p shifts from 0 to 1, there is some inversion point where the two are indifferent. Notice that we do not require global preferences \succeq to be complete. Nor do we require any form of independence which, together with the other axioms gives rise to the linear structure of expected utility, and has been controversial from the beginning. Allais' paradox, for instance, has been usually understood as an evidence of the failure of such independence axioms (cf. e.g. Rabin & Thaler (2001)).

The next postulate concerns social choice preferences.

A1 Convertibility. For every $x, y \in \mathcal{S}$, there exist $(z, \sigma) \in \mathcal{S} \times \text{Perm}(N)$ such that

$$z \sim_s x \text{ and } z\sigma \sim_s y.$$

A1 says that any pair of allocations in \mathcal{S} can be converted into an auxiliary pair of allocations related to each other by a permutation. Convertibility is implied by (but does *not* imply) the following selfishness assumption introduced by Kariv & Zame (2008):⁸

B1 selfishness. $x \sim_s (x_1, 0, \dots, 0)$ for every $x \in \mathcal{S}$.

Unfortunately, postulate **A1** on social choice preferences is not satisfied by many examples of interest (see subsection 4.2 *infra*). Therefore, we shall consider as well an alternative postulate on global preferences:

A2 Reduction of simple lotteries. For every simple lottery of the form $L = (px^j, (1-p)y^j) \in \mathcal{L}$, there exists a deterministic allocation $x \in \mathcal{S}$ such that:

$$x \sim \mathcal{L}.$$

A2 says that, in terms of preferences, every random allocation admits a certain equivalent. It is satisfied by most of the textbook preferences towards risk we are aware of (see, nevertheless, Example 4 *infra* for a counterexample).

THEOREM. 1) For all ethical preferences \succeq_e and for social preferences satisfying **A1**, there is a unique global preference relation \succeq on \mathcal{L} verifying **A0** such that its restriction to \mathcal{E} coincides with \succeq_e . Hence, if \succeq verify **A0** and are such that \succeq_s satisfy **A1**, then both risk preferences, \succeq_r , and social preferences, \succeq_s , can be deduced from ethical preferences \succeq_e .

2) For all ethical and social preferences, there is a unique global preference relation \succeq on \mathcal{L} verifying **A0** and **A2** whose restriction to \mathcal{E} coincides with \succeq_e . Hence, if \succeq verifies both **A0** and **A2**, then risk preferences and social preferences can be deduced from ethical preferences.

⁸To see that **B1** \Rightarrow **A1**, consider $z := (x_1, y_1, 0, \dots, 0)$ and $z\sigma := (y_1, x_1, 0, \dots, 0)$. **B1** implies that $z \sim_s x$ and $z\sigma \sim_s y$.

Proof. 1) Since $\mathcal{S} \subset \mathcal{E}$, social preferences can be deduced from ethical preferences. What we have to prove is that global preferences over \mathcal{L} can be deduced from ethical preferences \succeq_e (although, obviously, \mathcal{L} is not a subset of \mathcal{E}). Given assumption **A0**(i)-(iv) on global preferences \succeq , they verify the following property:⁹ For any lottery $(p_j x^j)_{j=1,\dots,K} \in \mathcal{L}$, it is possible to find a lottery involving only x^1 and x^K , and to which it is indifferent. That is, there exists $p \in [0, 1]$ such that

$$(p_j x^j)_j \sim (p x^1, (1-p)x^K).$$

Indeed, take $L := (p_j x^j)_{j=1,\dots,K} \in \mathcal{L}$. Continuity and substitutability imply that $L \sim (p_j X^j)_j$, where $X^j = (q_j x^1, (1-q_j)x^K)$ is defined as in (3). Reduction of compound lotteries and transitivity imply that $L \sim ((\sum_j p_j q_j)x^1, (\sum_j p_j(1-q_j))x^K)$. Therefore, for our purposes, it suffices to prove that the restriction of global preferences \succeq on simple lotteries of the form $(p x, (1-p)y)$ can be deduced from \succeq_e . Assumption **A1** enables us to find $(z, \sigma) \in \mathcal{S} \times \text{Perm}(N)$ such that $x \sim_s z$ and $y \sim_s z\sigma$. Given ethical preferences \succeq_e , let us therefore define global preferences by:

$$(p x^1, (1-p)y^1) \succeq (q x^2, (1-q)y^2) \iff (p z^1, (1-p)z^1\sigma_1) \succeq_e (q z^2, (1-q)z^2\sigma_2).$$

Clearly, global preferences defined this way will coincide with ethical preferences when restricted to \mathcal{E} , and hence with social preferences when restricted to \mathcal{S} . (Whenever $(p z^1, (1-p)z^1\sigma_1)$ and $(q z^2, (1-q)z^2\sigma_2)$ cannot be compared due to the lack of completeness of \succeq_e , then they cannot be compared either with respect to \succeq , which will therefore be incomplete.) On the other hand, since global (hence, also ethical) preferences are transitive, they are uniquely defined this way. Suppose, indeed, there are two pairs $(z, \sigma), (z', \sigma') \in \mathcal{E}$ such that

$$z \sim z' \sim x \text{ and } z\sigma \sim z'\sigma' \sim y.$$

Then, $z \succeq_e z\sigma \iff z' \succeq_e z'\sigma' \iff x \succeq y$. Therefore, the restriction of \succeq to \mathcal{R} yields a unique preference relation, \succeq_r , in the risk environment.

2) Take two simple lotteries $L_1, L_2 \in \mathcal{L}$ of the form $L_i = (p x^i, (1-p)y^i)_{i=1,2}$. By **A2**, there exist $x, y \in \mathcal{S}$ with $x \sim L_1$ and $y \sim L_2$. Define global preferences \succeq by

$$L_1 \succeq L_2 \iff x \succeq_e y.$$

□

Example 4 in subsection 4.2 will show that our Theorem is tight in the sense that one cannot relax both **A1** and **A2** without impairing our result.

4 Deducing ethics from economic decisions?

The previous section provided fairly weak assumptions under which risk and social preferences are uniquely determined by ethical ones. By contrast, we provide, now, a somewhat severe restriction that will be shown to imply the opposite property, that is, under which ethical preferences can be deduced from risk and social ones.

WI Weak independence. For every probability vector, (p_j) , and every pair of arrays of allocations $(x^j)_j$ and $(y^j)_j$, one has:

⁹See Luce and Raiffa (1957, p. 28).

$$x^j \succeq y^j \forall j \Rightarrow (p_j x^j)_j \succeq (p_j y^j)_j.$$

WI is a weakening of the familiar independence axiom, and does not imply expected utility (even combined with the rest of assumption **A0**).¹⁰

C Probabilistic self-regarding. 1) Let $(p_\sigma x^\sigma)$ and $(q_\sigma y^\sigma)$ be two lotteries in $\mathcal{E} \setminus \mathcal{S}$ such that $(p_\sigma x^\sigma) \succeq_e (q_\sigma y^\sigma)$. Then, there exists a pair, $(\tilde{x}^\sigma)_\sigma, (\tilde{y}^\sigma)_\sigma$, of allocations in \mathcal{S} such that: $(p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \sim (p_\sigma x^\sigma)$ and $(q_\sigma(\tilde{y}_1^\sigma, 0, \dots, 0)) \sim (q_\sigma y^\sigma)$. Moreover, for each such pair $(\tilde{x}^\sigma, \tilde{y}^\sigma)_\sigma$, one has:

$$(p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \succeq_r (q_\sigma(\tilde{y}_1^\sigma, 0, \dots, 0)).$$

Roughly speaking, condition **C** says that 1) every non-degenerate random allocation in \mathcal{E} is indifferent to some random allocation in \mathcal{R} , and 2) when evaluating a random allocation in \mathcal{E} , the Representative does not pay attention to the way randomness affects citizens different from herself. To put it differently, the attitude towards risk of citizens different from 1 has no impact on global preferences. We view this as a particularly severe restriction: How “ethical” are ethical preferences neglecting the risk aversion of the population’s vast majority ?

PROPOSITION 1.— If ethical preferences, \succeq_e , satisfy weak independence **WI**, then the two following conditions are equivalent:

- (a) \succeq_e can be deduced from \succeq_r and \succeq_s ;
- (b) \succeq_e verify **C**.

Proof.

(a) \Rightarrow (b). Suppose that (b) is not satisfied; we prove that (a) fails. Let us denote by $E \subset (\mathcal{E} \setminus \mathcal{S})$ the subset of random allocations $(p_\sigma x^\sigma)$ for which there exists $(\tilde{x}^\sigma) \in \mathcal{S}$ with $(p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \sim (p_\sigma x^\sigma)$. Suppose that $E \neq (\mathcal{E} \setminus \mathcal{S})$, and take $(p_\sigma x^\sigma)$ in $\mathcal{E} \setminus (E \cup \mathcal{S})$. Then, consider two global preferences, \succeq^1 and \succeq^2 , whose restrictions to \mathcal{R} (identified with a subset of \mathcal{L} thanks to Axiom SC) and \mathcal{S} both coincide with \succeq_r and \succeq_s , and such that, for some $(q_\sigma y^\sigma) \in \mathcal{E}$:

$$(p_\sigma x^\sigma) \succeq^1 (q_\sigma y^\sigma) \text{ while } (p_\sigma x^\sigma) \prec^2 (q_\sigma y^\sigma).$$

The restrictions to \mathcal{E} of \succeq^1 and \succeq^2 do not coincide although both global preferences are compatible with \succeq_r and \succeq_s . Hence \succeq_e is not uniquely determined by risk and social choice preferences.

Suppose, next, that $E = \mathcal{E}$ but there exists a pair of lotteries, $(p_\sigma x^\sigma), (q_\sigma y^\sigma)$, with the property that $(p_\sigma x^\sigma) \succeq (q_\sigma y^\sigma)$ and yet $(p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \prec_r (q_\sigma(\tilde{y}_1^\sigma, 0, \dots, 0))$ for some allocations $(\tilde{x}^\sigma, \tilde{y}^\sigma) \in \mathcal{S}$ verifying $(p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \sim (p_\sigma x^\sigma)$ and $(q_\sigma(\tilde{y}_1^\sigma, 0, \dots, 0)) \sim (q_\sigma y^\sigma)$.¹¹ Consider the ethical preference, \succeq_e^* , defined by:

¹⁰Weak independence is sometimes known as the “sure thing principle”, and is essentially identical to the game-theoretical principle that a rational individual will avoid using any weakly dominated strategy.

¹¹Recall that \succeq need not be transitive.

$$(p_\sigma x^\sigma) \succeq_e^* (q_\sigma y^\sigma) \iff (p_\sigma(\tilde{x}_1^\sigma, 0, \dots, 0)) \succeq_r (q_\sigma(\tilde{y}_1^\sigma, 0, \dots, 0)).$$

Then, $\succeq_e^* \neq \succeq_e$ although they are both compatible with \succeq_r and \succeq_s . Hence, given risk and social preferences, ethical preferences are not uniquely defined.

(b) \Rightarrow (a). It suffices to define \succeq_e^* as above, and to conclude from weak independence that ethical preferences are uniquely defined once \succeq_r and \succeq_s are given. \square

Let us now recall the restrictions on social preferences introduced by Kariv and Zame (2008). In addition to being complete, transitive, reflexive, and continuous, they need to verify:

B2 The worst outcome. $x \succeq_s 0$ for every $x \in \mathcal{S}$.

This requirement is specific to their framework as they impose allocations to take value in \mathbb{R}_+^N . No such restriction is needed in our set-up.

B3 Self-regarding. For each $x \in \mathcal{S}$, there is a $t \in \mathbb{R}_+$ such that $(t, 0, \dots, 0) \succeq_s x$.

Clearly, selfishness **B1** is a strengthening of “self-regarding” **B3**. The two results proven in Kariv and Zame (2008) that are of interest to us are the following:

PROPOSITION 2.— (Kariv and Zame (2008)) 1) For all risk preferences and social preferences that satisfy **B2** and **B3**, there is a unique preference relation \succeq on \mathcal{L} verifying **WI** and whose restriction to \mathcal{S} (resp. \mathcal{R}) coincides with \succeq_s (resp. \succeq_r). Hence, if \succeq verifies Weak independence, then ethical preferences can be deduced from risk and social preferences.

2) If social preferences are selfish (i.e., satisfy **B1**), then \succeq has the following property: For all lotteries $(p_j x^j), (q_k y^k) \in \mathcal{L}$,

$$(p_j x^j)_j \succeq (q_k y^k)_k \iff (p_j(x_1^j, 0, \dots, 0))_j \succeq_r (q_k(y_1^k, 0, \dots, 0))_k.$$

It is easily shown that, if **WI**, **B2** and **B3** are fulfilled, so is **C**. Hence, Part 1 of Proposition 2 follows from Part 2 of our Proposition 1. The second part of Proposition 2 says that, if the Representative is selfish (in the sense of **B1**) in the social choice environment, then preferences in the risk environment coincide with ethical preferences. Given the widespread use of expected utility, this seems *prima facie* to promote a definition of ethical preferences as being given by the expected utility of random allocations x^σ for $\sigma \in \text{Perm}(N)$:

$$U[(p_\sigma x^\sigma)] := \sum_{\sigma} p_{\sigma} \sum_i \lambda_i x_{\sigma(i)}. \quad (4)$$

Notice, however, that Proposition 2 is hardly compatible with (4) since this criterion does *not* verify selfishness **B1** unless the weights $(\lambda_i)_i$ attributed to citizens are $\lambda_1 > 0$ and $\lambda_i = 0$ for every $i \neq 1$ —in which case (4) simply reduces to dictatorship ! Nor would (4) fulfill **B2** once prospects are allowed to take values that are unbounded from below. Whether (4) can be understood in terms of Harsanyi’s “utilitarian ethics” is discussed below in Example 2.

5 Examples

The first three examples satisfy our axioms but fail to verify **C** (hence the assumptions adopted by Kariv and Zame (2008) as well). The last example shows that our main result fails if Axiom **A** is abandoned.

Example 1. The Maximin criterion (both with respect to risk and with respect to citizens) can be defined by the global utility function on \mathcal{L} :

$$M((p_k x^k)_k) := \min_{k,i} x_i^k. \quad (5)$$

Axioms **C**, **A0** are verified but neither **B2**, nor **B3** (hence **B1**), and **C**.¹² Moreover, even under the Strong Compatibility assumption, ethical preferences \succeq_e cannot be deduced from risk preferences \succeq_s when allocations are constrained to take values in \mathbb{R}_+^N . Indeed, the restriction of (5) to allocations of the form $(x_1, 0, \dots, 0)$ yields a constant mapping, so that risk preferences are trivial. On the other hand, consider the auxiliary utility function on global lotteries:

$$N((p_k x^k)_k) := \min \left\{ \min_k x_1^k ; \min_{i \neq 1} \sum_k p^k x_i^k \right\}. \quad (6)$$

The restriction of (6) to \mathcal{R} and \mathcal{S} yields the same risk and social preferences as (5), while the restrictions of both global utilities to the ethical environment, \mathcal{E} , are distinct. Hence, ethical preferences cannot be deduced from \succeq_r and \succeq_s .

Following Harsanyi (1975, 1978), one has argued that the risk preferences induced by (5) in the \mathcal{R} setting are hardly realistic. Quite on the contrary, both theoretical investigations (see, e.g., Artzner *et al.* (1999)) and empirical practices of stress tests in the financial industry suggest that behaviors at least close to the ones dictated by (5) are not relegated to exotic matters, even in the highly specific set-up of individual risk. Similarly, Gilboa and Schmeidler (1989) have reintroduced the maximin principle within decision theory in face of uncertainty. On the (deterministic) social choice side, such an egalitarian criterion has been strongly advocated by Fleurbaey and Maniquet (2006) in a purely ordinal setting.

Finally, we must admit that it is unclear whether Rawls would agree with our use of probabilities (which are absent from his description of the “original position”). We do not lay claim to Rawlsian orthodoxy on this point.

Example 2. Consider lotteries involving at most $K \geq 2$ allocations,¹³ $(x^k)_k$, ordered so that $p_{k+1} \leq p_k$ for $k = 1, \dots, K-1$. A criterion akin to some kind of “utilitarianism” (again, both with respect to risk and to citizens) can be defined by:

¹²This would be true also for Leximin preferences as well.

¹³We know from the proof of the Theorem that this involves no loss of generality.

$$U((p_k x^k)_k) := \sum_k p_k \sum_i \lambda_i^k x_i^k, \quad (7)$$

with $\lambda_i^k \in \mathbb{R}$. This alternate criterion fulfills **C**, **A0** and **A2** but fails to verify **B3**, **B1** and **A1**. Moreover, when the individual weights λ_i^k depend upon k in a non-trivial way, **C** is not satisfied either, so that ethical preferences \succeq_e cannot be deduced from risk \succeq_r and social \succeq_s preferences. Indeed, neither \succeq_r nor \succeq_s depend upon λ_i^k for $k \geq 2$ and $i \geq 2$. On the side of individual risk, (7) corresponds to risk-neutrality which is widely used for pricing and hedging financial derivatives. On the side of social choice, it has received an axiomatic foundation by Mertens and Dhillon (1999).

There has been considerable controversy over “utilitarian ethics” in the way it is defended by Harsanyi, as in the debate between Sen (1976, 1977, 1986) and Harsanyi (1975, 1977a). Here, when $p_\sigma = 1/N!$ for each permutation $\sigma \in \text{Perm}(N)$, our framework becomes compatible with Harsanyi’s (1975) “equi-probability model of moral value judgments”. To see this point, recall that, in Harsanyi’s (1978) view, the Representative “would certainly satisfy our impartiality and impersonality requirements if he did not know how his choice between [lotteries] A and B would affect him personally and, in particular, if he did not know what his own social position would be in situations A and B ”. Thus, the Representative is assumed to think that in either (randomly selected) situation he would have the same probability $1/N$ to occupy any one of the N possible social positions. Therefore, in Harsanyi (1978), the Representative does not even know her own risk preferences, as these preferences are attached to the position she will occupy, while, here, a lottery in \mathcal{E} involves various random prospects to individuals $i = 1, \dots, N$ (including the Representative) knowing her own, fixed, risk preferences, \succeq_r^i .

Nevertheless, our approach is broad enough to encompass Harsanyi’s set-up as a particular case of ours: Suppose that the index $i = 1, \dots, N$ does not label individuals but “social positions”, which may be occupied by every individual. Take $L = 1$ and suppose that each position $i = 1, \dots, N$ is identified with a given utility function: $U_i : \mathcal{A} \rightarrow \mathbb{R}$ defined on some auxiliary space, \mathcal{A} , of random situations, $(p_k A^k)_k$. An allocation of prospects, $(x_i)_i \in \mathbb{R}^N$, is now a N -tuple of utility levels $(U_i(A))_i$, derived from any random situation $A \in \mathcal{A}$.¹⁴ Restrict \mathcal{E} to equiprobable lotteries, $(p_\sigma x^\sigma)_{\sigma \in \text{Perm}(N)}$, of size $N!$, with $p_\sigma = 1/N!$, every σ . By construction, a form of “selfishness” is implicit in Harsanyi’s framework since, whatever being her position i , the Representative only cares about her own individual risk preferences, U_i , associated to this very position, and not about the preferences of the other citizens occupying different positions — so that **B3**, now, makes sense. Within this specific set-up, Proposition 2-2 provides a first step towards Harsanyi’s conclusion. It suffices, indeed, to complete the assumptions needed for Proposition 2-2 by any axiomatics which characterizes individual risk preferences in terms of expected utility

¹⁴ Admittedly, this construction involves interpersonal utility comparisons — which is consistent with Harsanyi’s (1977b, 1978) claim that “there are non valid arguments against such comparisons”. Though we do believe that there are valid arguments against intersubjective comparisons (whose discussion would go beyond the scope of this paper and is abundantly illustrated in the literature), it is only fair to permit them in order to characterize Harsanyi’s setting as a particular instance of ours.

in order to get:

$$\begin{aligned} U_H[(p_\sigma x^\sigma)] &:= \frac{1}{N!} \sum_{\sigma} x_{\sigma(1)} \\ &= \frac{1}{N} \sum_i x_i \\ &= \frac{1}{N} \sum_i U_i(A). \end{aligned}$$

The advantage of this reformulation is to illuminate the role of the selfishness requirement **B1** underlying this “utilitarian”¹⁵ approach of ethics.¹⁶

Example 3. (Kariv and Zame (2008)) Take $N = 2$. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be any continuous, strictly increasing function with the property that $g(t) = t$ for $t \leq 0$. Define the global utility function $W_g : \mathcal{L} \rightarrow \mathbb{R}$ by:

$$W_g(p_1(x_1, y_1), p_2(x_2, y_2)) := p_1 g(-e^{x_1} + y_1) + p_2 g(-e^{x_2} + y_2),$$

for any simple lottery in \mathcal{L} involving only two allocations (x_1, y_1) and (x_2, y_2) with $p_1 \geq p_2$.¹⁷ The restriction of W_g on \mathcal{R} does not depend on g since $g(t) = t$ for $t \leq 0$. The social preferences induced by W_g on \mathcal{S} do not depend upon g because g is strictly increasing. However, the ethical preferences induced by W_g on \mathcal{E} do depend on g : The weight given to inequality between citizens depends on g . Hence, ethical preferences *cannot* be deduced from risk and social choice preferences, so that Proposition 1 above fails. This is due to the failure of **B3**: Preferences induced by W_g on \mathcal{S} are not self-regarding. Neither are they probabilistically self-regarding, so that **C** is, in turn, violated. By contrast, the intermediate value theorem ensures that W_g verifies **A2**, while **A0** is obvious. Hence, our Theorem holds in this setting.

Example 4. Again, take $N = 2$. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be continuous, strictly increasing with the property that $f(t) = t$ whenever $t \geq 0$. Define the global utility function $U_{f,\lambda} : \mathcal{L} \rightarrow \mathbb{R}$ as follows. For every lottery $L = (p_1(x_1, y_1), p_2(x_2, y_2))$ with $p_1 \geq p_2$:

$$U_{f,\lambda}(L) := p_1 f(|x_1| - y_1) + p_2 \lambda f\left(\max_{i=1,2} |x_i - y_i| - \min_{i=1,2} |x_i - y_i|\right),$$

for a given parameter $\lambda \in \mathbb{R} \setminus \{0\}$. The risk preferences induced by $U_{f,\lambda}$ on \mathcal{R} reduce to

$$U_{f,\lambda}(L) = p_1 f(|x_1|) + p_2 \lambda f\left(\max_{i=1,2} |x_i| - \min_{i=1,2} |x_i|\right),$$

for every $L := (p_1(x_1, 0), p_2(x_2, 0))$ and do not depend on f (as $f(t) = t$ for $t \geq 0$). Nor do social choice preferences since, for every allocation $(x_1, y_1) \in \mathcal{S}$, $U_{f,\lambda}(x_1, y_1) = f(|x_1| - y_1)$, and f is strictly increasing. However, the ethical preferences induced

¹⁵Quotation marks, here, wish to emphasize that Harsanyi’s terminology does not reflect the much broader standpoint of, say, John Stuart Mill (1861), for whom “utilitarianism” also encompasses non-self-oriented behaviors (e.g., the Biblical Golden rule) which obviously contradict **B3**.

¹⁶See the discussion of (4) *supra*.

¹⁷By the same argument as in the proof of our Theorem, it suffices to consider such simple lotteries.

by $U_{f,\lambda}$ on \mathcal{E} do depend on f , again because the weight given to inequality depends upon f . Thus, ethical preferences cannot be deduced from risk and social choice preferences. This time, it is **B2** that fails: There is no worst outcome. On the other hand, that **C** is not satisfied is obvious. At variance with Example 3, however, $U_{f,\lambda}$ does *not* verify our convertibility assumption **A1**. Similarly, $U_{f,\lambda}$ does not satisfy **A2** in general, so that our Theorem fails as well. Indeed, for every allocation $z = (x, y)$, $U_{f,\lambda}(pz, (1-p)z\sigma) = pf(|x| - y)$ whatever being the probability $p \in (0, 1]$ with $p \geq 1 - p$.¹⁸ Hence, ethical preferences do not depend upon λ , while global preferences do. In particular, risk preferences depend upon λ , and hence, cannot be deduced from ethical ones.

6 Discussion

What do we concretely *mean* when we claim that economic decisions can be “deduced” from ethical ones? Obviously, we do not intend to provide a substantialist definition of ethics in the Aristotelian style. Rather, this paper confines itself within the Kantian (liberal) tradition according to which every person is responsible for deciding what (ethical) “good” means for her. Nevertheless, following Rawls, Harsanyi and many other social philosophers (such as, e.g., J. Habermas (1991)), we view ethical decisions behind the veil of ignorance as being potentially discussed in the public space of debates, on the basis of rational arguments such as those encapsulated in the various axiomatics of social choice theory already alluded to.¹⁹ Claiming that economic decisions depend upon ethical ones therefore concretely means that economic and social decisions should be included in the agenda of ethical public debates. They cannot be derived from individual risk preferences, from which they would inherit the undebatable arbitrariness. If one adopts the axioms of Theorem 1, then choosing among various political and economic systems is *not* the same kind of decision as choosing among various financial portfolios or various menus in a restaurant. If, now, individual risk preferences, \succeq_r , are interpreted as preferences of a politician who faces various voting systems and aims at maximizing the chance of being (re)elected, this paper also shows that political decisions cannot be *reduced* to the self-interested calculus of politicians, contrary to what tends to assume the public choice school (cf. e.g., Buchanan & Tullock (1962)). This does not mean that we deny any relevance to the public choice approach: All we claim is that such a self-interested calculus cannot exhaust the freedom of politicians and cannot exonerate their decisions from any ethical questioning — unless one is ready to adopt the restrictive assumptions of subsection 4.1.

How “realistic” is our conclusion that ethical preferences cannot be deduced from every-day behavior on the economic field? Many contemporary “utilitarians” have claimed that voting for the maximin principle is only optimal for infinitely risk averse Representatives. This argument takes as granted that each person is only interested in her own material payoff — not surprisingly, this is exactly assumption **B1** — and claims that it is legitimate to disregard the Maximin principle on the ground that people’s everyday behavior does not fit with infinite risk-aversion. This presupposes

¹⁸Here, σ denotes the unique non-trivial permutation over $\{1, 2\}$.

¹⁹See, in particular, Mertens & Dhillon (1999) and Fleurbaey & Maniquet (2008).

exactly what this paper challenges, namely that ethical preferences can be deduced from risk preferences.

Hörisch (2007) has implemented the Rawlsian thought experiment of a veil of ignorance as a laboratory experiment. There, it is found that both men and women react to the risk introduced by the veil of ignorance in a way that is significantly distinct from their attitude towards risk in front of the veil. Women additionally exhibit social preferences that reflect an increased concern for equality. These findings confirm the main message of the present paper. Indeed, if people have social preferences that do not satisfy **B3**, they can be in favor of an egalitarian distribution even if they are, say, risk-neutral.

Conversely, one could question the “realism” of our assertion that economic decisions are influenced by ethical convictions. In a sense, the main message of this paper goes in the same direction as the one put on the forefront, e.g., by Foley (2008): Economic theory at its most abstract level is a speculative philosophical discourse, shaped by (at least implicit) ethical options as well as by deductive or inductive scientific findings. The attempt to separate the economic sphere of life, in which the pursuit of self-interest is led by some invisible hand of the market to a socially beneficial outcome, from the rest of social life, in which the pursuit of self-interest is morally problematic and has to be weighed against other ends, this attempt is a short-cut that, according to the results given in this paper, need to be reversed: the economic sphere of life is part of social life.

Demichelis & Weibull (2008) have proven that “honesty” enables evolutionary stable Nash equilibria of (the lexicographic communication extension of) generic and symmetric $n \times n$ -coordination games to concentrate on the unique Pareto optimal outcome of the underlying game. This shows, at least, that there exist situations where ethics matters even at the efficiency level. In fact, ethical views *de facto* influence investor as well as consumer choices, not only at the time of presbyterian pietism studied in Weber’s (1904) celebrated monograph, but also today. For example, Freeman (2001) argues that Rawls’ original position could be used as a valuable thought experiment for orienting management decisions. On the other hand, in 2006, the UN launched an initiative called “Principles for Responsible Investments”:²⁰ The asset owners and investment managers who sign the six principles commit themselves to integrate ESG (environmental, social and governance) criteria in their investment decisions. By May 2008, 362 investors had signed these principles, representing 14.4 trillion dollars of investments. Fair Trade is also an alternative way of doing business that seeks to build equitable, long-term partnerships between consumers in Europe, Japan, Australia, New Zealand and North America together with producers in developing regions. The global Fair Trade sales in 2007 are worth 2.65 billion euros..²¹

Finally, current initiatives in favour of social business (Yunus (2008)) express the will of entrepreneurs to endorse new economic models centered on the needs of the poor, even if these actions are less profitable than conventional businesses.

Let us conclude with a final remark. Nussbaum (2006, p. 17) criticizes the social contract theorists, and Rawls among them, in as much as they see the society as a contract for mutual advantage between people who are free, equal and independent.

²⁰See <http://www.unpri.org>.

²¹Cf. Krier (2007).

This perspective does not take into account people who suffer from impairments or disabilities. Even though we agree with Nussbaum's criticism, we did not address the issue in this paper: The parties behind the veil of ignorance do not possess any serious physical or mental impairments that would prevent them from exhibiting "preference" relations fulfilling **C** and **A0** together with either **A1** or **A2**. However, the citizens for whom they design principles could suffer from such disabilities.

GAËL GIRAUD, CNRS, PARIS SCHOOL OF ECONOMICS, ESCP-Europe, gael.giraud@parisschoolofeconomics.eu.
CÉCILE RENOUARD, ESSEC BUSINESS SCHOOL, PARIS, renouard@essec.fr.

References

- [1] Artzner, Ph., Fr. Delbaen, J.-M. Eber, D. Heath "Coherent Measures of Risk", *Mathematical Finance* 9 n. 3 (1999) 203-228
- [2] Buchanan, J. M. & G. Tullock (1962) *The Calculus of Consent: Logical Foundations of Constitutional Democracy*, Ann. Arbor, University of Michigan Press.
- [3] Clark, A.E. & A.J. Oswald (1996) "Satisfaction and Comparison Income", *Journ. of Pub. Econ.*, 61, 359-381.
- [4] Demichelis, S. & J. Weibull (2008) "Language, Meaning, and Games: A Model of Communication, Coordination, and Evolution", *American Economic Review*, 98:4, 1292-1311.
- [5] Easterlin, R. A. (1974) "Does Economic Growth Improve the Human Lot?" in Paul A. David and Melvin W. Reder, eds., *Nations and Households in Economic Growth: Essays in Honor of Moses Abramovitz*, New York: Academic Press, Inc.
- [6] Fehr, E. and K. Schmidt (2006) "The Economics of Fairness, Reciprocity and Altruism — Experimental Evidence and New Theories", in Kolm, S.-C., Ythier, J.-M. (Eds) *Handbook of the Economics of Giving, Reciprocity and Altruism*, 2006, vol. I. forthcoming.
- [7] Fleurbaey, M. and Fr. Maniquet (2008) "Fair Social Orderings", *Economic Theory*, 34, 25-45.
- [8] Foley, D. (2008) *Adam's Fallacy, a Guide to Economic Theology*, Belknap press.
- [9] Freeman, R. E. (1994), "The Politics of Stakeholder Theory: some Future Directions", *Business Ethics Quarterly*, volume 4(4), 409-421.
- [10] Gilboa, I. and D. Schmeidler "Maximin Expected Utility with Non-unique Prior", *Journ. of Mathematical Economics*, 18 (1989) 141-153.
- [11] Habermas, J. (1991) *Moral Consciousness and Communicative Action*, Cambridge: MIT Press.
- [12] Harsanyi, J. "Cardinal Utility in Welfare Economics and in the Theory of Risk-taking", *Journal of Political Economy*, 61 (1953) 434-435.
- [13] ——— "Cardinal Welfare, Individualistic Ethics, and Interpersonal Comparisons of Utility", *Journal of Political Economy*, 63 (1955) 309-321.

- [14] ——— “Can the Maximin Principle Serve as a Basis for Morality ? A Critique of John Rawls’ Theory”, *American Political Science Review*, 69 (1975) 594-606.
- [15] ——— “Nonlinear Social Welfare Function: A Rejoinder to Professor Sen”, in *Foundational Problems in the Social Sciences*, R. Butts and J. Hintikka (eds.), 1977a, 293-296.
- [16] ——— “Morality and the Theory of Rational Behavior”, *Social Research*, 44 (1977b) 623-656.
- [17] ——— “Bayesian Decision Theory and Utilitarian Ethics”, *American Economic Review*, 68 (1978) 223-228.
- [18] Hörisch, H. “Is the Veil of Ignorance only a Concept about Risk? An Experiment”, University of Munich (2007) (discussion paper).
- [19] Kariv, S. and W. Zame “Piercing the Veil of Ignorance”, mimeo (2008) <http://emlab.berkeley.edu/~kariv/KZ.I.pdf>
- [20] Krier, J.-M. “Fair Trade 2007: New Facts and Figures from an Ongoing Success Story. A Report on Fair Trade in 33 Consumer Countries” (2007) <http://www.european-fair-trade-association.org/efta/Doc/FT-E-2007.pdf>
- [21] Mertens, J.-F. and A. Dhillon (1999) “Relative Utilitarianism”, *Econometrica*, 3, 471-498.
- [22] Mill, J. S. *Utilitarianism*, Every’s Man Library, [1861] 1992, London.
- [23] Nussbaum, M. *Frontiers of Justice*, Harvard University Press, 2006.
- [24] Rabin, M. & R. H. Thaler (2001) “Risk Aversion,” *Journal of Economic Perspectives*, Vol. 15, N.1, 219-232.
- [25] Rawls, J. *A Theory of Justice*, Cambridge, Harvard University Press, 1971.
- [26] ——— “Reply to Alexander and Musgrave”, *Quarterly Journal of Economics*, 88 (4) (1974) 633-655.
- [27] Ricœur, P. *Oneself as Another*, transl. K. Blamey, Chicago, University of Chicago Press, 1992.
- [28] Segal, U. (2000). ‘Lets Agree That All Dictatorships are Equally Bad,’ *J. Polit. Econ.* 108, 56989.
- [29] Sen, A. K. “Welfare Inequalities and Rawlsian Axiomatics”, *Theory and Decision*, 7 (1976) 243-262.
- [30] ——— “Non-linear Social Welfare Functions: A Reply to Professor Harsanyi”, in *Foundational Problems in Social Sciences*, R. Butts and J. Hintikka (eds.), 1977, 297-302.
- [31] ——— “Social Choice Theory”, in *The Handbook of Mathematical Economics*, vol. III, K. Arrow and M. Intrilligator (eds.), 1986, 1073-1181.
- [32] Sobel, J. (1981). “Distortion of Utilities and the Bargaining Problem,” *Econometrica* 49, 597619.
- [33] ——— (2001) “Manipulation of Preferences and Relative Utilitarianism” *Games and Economic Behavior* 37, 196215.

- [34] Weber, M. *The Protestant Ethic and the Spirit of Capitalism*, Penguin Books, [1904] (2002) transl. P. Baehr and G. C. Wells.
- [35] Yunus M. *Creating a World without Poverty: Social Business and the Future of Capitalism*, 2008, Public Affair Pr.