



WARWICK

ECONOMICS

CRETA

Centre for Research in Economic Theory and its Applications

Discussion Paper Series

**Fairness and Utilitarianism without Independence**

Sinong Ma & Zvi Safra

March 2016

No: 20

CRETA

Centre for Research in Economic Theory and its Applications

Department of Economics  
University of Warwick, Coventry,  
CV4 7AL, United Kingdom

[warwick.ac.uk/fac/soc/economics/research/centres/creta/papers](http://warwick.ac.uk/fac/soc/economics/research/centres/creta/papers)

# Fairness and Utilitarianism without Independence

Sinong Ma<sup>1 2</sup>

Zvi Safra<sup>3 4</sup>

March 6, 2016

<sup>1</sup>The order of authors names was determined by a toss of a fair coin...

<sup>2</sup>University of Warwick; email: phd12sm@mail.wbs.ac.uk.

<sup>3</sup>University of Warwick (Emeritus, Tel Aviv University); email: z.safra@wbs.ac.uk.

<sup>4</sup>We are grateful to Simon Grant, Edi Karni, Tigran Melkonyan, Uzi Segal and Jean-Marc Tallon for very useful comments and suggestions.

## **Abstract**

In this work we reconsider Harsanyi's celebrated (1953, 1955, 1977) utilitarian impartial observer theorem. Departing from Harsanyi's individual-centered approach, we argue that, when societal decisions are at stake, postulates must not be drawn from individualistic behavior. Rather, they should be based on societal norms. Hence, notions like societal fairness should explicitly be taken as the guiding principles. Continuing this line of thinking, we state and prove a utilitarian result that, rather than the independence assumption, is based on the notion of procedural fairness and on similar treatment of societal and individual lotteries.

“An axiomatic justification of utilitarianism would have more content to it if it started off at a place somewhat more distinct from the ultimate destination” Sen (1976, page 251)

## 1 Introduction

In this work we reconsider Harsanyi’s celebrated (1953, 1955, 1977) utilitarian impartial observer theorem. We propose an approach that puts more emphasis on procedural fairness and offer a utilitarian result that does not use the independence assumption.

Harsanyi analyzed a society that needs to choose among alternate social policies, each of which is a probability distribution (a lottery  $\ell$ ) over a given set of social actions, where the latter associate outcomes with the society’s members. Every social lottery  $\ell$  induces a lottery  $\ell_i$  on individual  $i$ . Individual  $i$ ’s preferences  $\succsim_i$  are known and different individuals may possess distinct preferences.

To help determine the optimal social policy, Harsanyi suggested that every individual is endowed with social preferences. Individuals may develop these preferences by adopting the role of an impartial observer, thus disregarding their true identities and acting behind “a veil of ignorance”. Therefore, the impartial observer faces not only a lottery  $\ell$  over social actions, but also a lottery  $\alpha$  over identities. By assumption, the impartial observer is able to compare situations in which two individuals get two different outcomes, under two disjoint social actions.

Harsanyi argued strongly for “Bayesian rationality”. That is, he assumed that all individuals satisfy the *independence assumption* of the expected utility theory, both at their personal and social preference layers. Harsanyi claimed that these “sound” axioms, together with the so-called *acceptance principle* (that an impartial observer fully adopts individual  $i$ ’s preferences if she imagines becoming that individual for sure), would force the impartial observer to be a (weighted) utilitarian. More formally, over all extended lotteries  $(\alpha, \ell)$  in which the identity lotteries and action lotteries are independently distributed, the impartial observer’s preferences admit the following representation:

$$V(\alpha, \ell) = \sum_{i \in \mathcal{I}} \alpha_i U_i(\ell_i)$$

where  $\alpha_i$  is the probability of assuming person  $i$ 's identity and  $U_i(\ell_i) := \sum_x u_i(x)\ell_i(x)$  is person  $i$ 's von Neumann-Morgenstern expected utility.

Like Harsanyi, most authors who derived modifications of the utilitarianism result always assumed the independence axiom. See the works of Weymark (1991), Zhou (1997), Karni (1998), Dhillon and Mertens (1999), Gilboa, Samet and Schmeidler (2004), Grant, Kajii, Polak and Safra (2010; henceforth GKPS), Fleurbaey and Mongin (2012) and others. Notable exceptions Blackorby, Donaldson and Mongin (2004) and Mongin and Pivato (2015), who derived utilitarianism without independence. However, unlike the works mentioned above (including ours), these authors consider both *ex post* and *ex ante* analyses (and thus are able to employ Gorman's (1968) separability theorem).

Interestingly, Harsanyi's entire emphasis on Bayesian rationality was based on an individual-centered approach. Firstly, he assumed that rational individuals must satisfy the independence assumption and secondly, he claimed that society, by its need to be at least as rational as its individuals, must also satisfy independence (Harsanyi 1975). We disagree with Harsanyi on this. Instead we argue that when societal decision problems are at stake, postulates must not be drawn from individualistic behavior. Rather, they should be based on societal norms. Hence, when social preferences are formed, issues like societal fairness and equity should explicitly be taken to be the guiding principles.

In this work we focus on procedural fairness. This principle was first advocated by Diamond (1967) and was strongly supported by Sen (e.g., 1977). Its essence can be illustrated by the following example, which is an adoption of Diamond's example to the impartial observer framework. Consider a society that needs to decide on how to allocate an indivisible good among two individuals, 1 and 2, and action  $a^i$  denotes allocating it to individual  $i$ . Suppose, as Diamond did, that  $u_i(a^i) = 1$  for both  $i$ ,  $u_i(a^j) = 0$  for  $i \neq j$ , and that the impartial observer evaluations of all four outcomes are identical to those of the relevant individuals. The impartial observer imagines having equal chance of being either individual and has two policies at hand: policy (1) allocates the good to individual 1 (that is, chooses action  $a^1$ ) and policy (2) allocates the good to either individual by tossing a fair coin (that

is, choose the action lottery  $\frac{1}{2}a^1 + \frac{1}{2}a^2$ ). Clearly, Harsanyi’s utilitarian observer is indifferent between these two policies, as the expected social welfare values are identical. However, Diamond, and Sen, argued that the impartial observer might prefer policy (2), as it provides both individuals with a “fair shake”.<sup>1</sup> This notion of procedural fairness is expressed in our work through the notion of *convexity* over action lotteries: if, given an identity lottery  $\alpha$ , two individuals disagree on the ranking of action lotteries  $\ell$  and  $\ell'$ , then mixtures of these lotteries are weakly preferred over the less favorable one. If, as is the case in Diamond’s example, the impartial observer is indifferent between  $\ell$  and  $\ell'$ , then their mixture is weakly preferred to both.<sup>2</sup>

Working in a framework in which the basic building blocks are two different types of lotteries, those over identities and those over actions, raises a natural question: should these separate types be treated similarly? Harsanyi, by applying the independence axiom to all possible pairs of identity-action lotteries, implicitly assumed that they should. As we agree with Harsanyi on this, we make this assumption explicit and call it *indifference between identity and action lotteries*.<sup>3</sup>

Despite its innocuous appearance, the conjunction of this assumption with our notion of procedural fairness turns out to be rather forceful. More precisely, the main result of this work shows that convexity, indifference between identity and action lotteries and a stronger notion of acceptance are necessary, and sufficient, for utilitarianism.

Since the independence axiom is not assumed here, even not partially, this result is novel and quite unexpected. Paraphrasing Sen’s quote, we believe that one could hardly find an axiomatic justification of utilitarianism that starts off at a place that is more distinct from the ultimate destination than ours.

Lastly, our result implies that indifference between identity and action lotteries cannot hold if societies wish to exhibit *strict* inclination towards procedural fairness. Therefore, to accommodate views of authors like Diamond and Sen, the impartial observer must prefer

---

<sup>1</sup>A long list of real-life applications supporting Diamond’s fairness consideration is provided by Elster (1989).

<sup>2</sup>Unlike Epstein and Segal (1992), we do not assume strict preference because, as was argued by Sen (1977), mixture is not always superior.

<sup>3</sup>This assumption is similar to, though weaker than, an assumption made by GKPS (2010).

action lotteries. We elaborate on this in the concluding section.

This work is organized as follows: Section 2 sets up the framework, Section 3 presents the assumptions, Section 4 states, and explains, the utilitarian result and section 5 concludes. Finally, proofs are given in the appendix.

## 2 Setup and Notation

Let  $\mathcal{X} = [x_{\min}, x_{\max}] \subset \mathbb{R}$  be a compact interval representing all possible outcomes and let  $\Delta(\mathcal{X})$  denote the set of outcome lotteries, endowed with the weak convergence topology. With slight abuse of notation, we will let  $x$  denote the degenerate outcome lottery that assigns probability 1 to outcome  $x$ . Let  $T$  be a finite set of potential individual types, where each type  $t \in T$  is characterized by a preference relation  $\succsim_t$  defined over  $\Delta(\mathcal{X})$ . We assume throughout that each  $\succsim_t$  is complete, transitive, continuous (in that the weak upper and weak lower contour sets are closed in the product topology), its asymmetric part  $\succ_t$  is nonempty and is increasing with respect to first-order stochastic-dominance. The set of individuals under consideration is  $\mathcal{I} = \cup_{t \in T} \mathcal{I}_t$ , where  $\mathcal{I}_t$  is a denumerable (infinite) set of type  $t$  individuals. A society  $I$  is a finite subset of  $\mathcal{I}$ . Note that, even though we allow for societies in which all individuals are of the same type, these individuals need not be treated similarly. Also note that our framework departs from Harsanyi's in that, instead of working with one fixed finite society, we consider all societies that are subsets of  $\mathcal{I}$ .<sup>4</sup>

A social policy, or an action, associates an outcome with every individual and hence is represented by a function  $a : \mathcal{I} \rightarrow \mathcal{X}$ . The set of all these actions, endowed with the corresponding product topology, is denoted by  $\mathcal{A}$  (two extreme actions,  $a_{\max}$  and  $a_{\min}$ , defined by  $a_{\max}(i) = x_{\max}$  and  $a_{\min}(i) = x_{\min}$  for all  $i$ , respectively, will be used in the sequel). Let  $\Delta(\mathcal{A})$  denote the set of simple lotteries (lotteries with finite support) over actions, with typical elements denoted by  $\ell$ . With slight abuse of notation, we will let  $a$  denote the degenerate action lottery that assigns probability 1 to action  $a$ . A lottery  $\ell \in \Delta(\mathcal{A})$  is sometimes written as  $\ell = \sum_{a \in \text{Supp}(\ell)} \ell(a) a$ .

Following Harsanyi, an observer imagines herself behind a veil of ignorance, uncertain

---

<sup>4</sup>The need for an infinite set of individuals is clarified in the proof of the theorem.

about which identity she will assume in the given society.<sup>5</sup> Let  $\Delta(\mathcal{I})$  denote the set of simple identity lotteries on  $\mathcal{I}$ , where typical elements are denoted by  $\alpha$  ( $\alpha_i$  denotes the probability assigned by the identity lottery  $\alpha$  to individual  $i$ ). These lotteries represent the imaginary risks in the mind of the observer of being born as someone else. With slight abuse of notation, we will let  $i$  denote the degenerate identity lottery that assigns probability 1 to individual  $i$ . An imaginary lottery  $\alpha \in \Delta(\mathcal{I})$  is sometimes written as  $\alpha = \sum_{i \in \text{Supp}(\alpha)} \alpha_i i$ . When the observer is faced with pairs of identity and action lotteries, it is assumed that they are independently distributed.

The observer is endowed with a preference relation  $\succsim$  defined over the space of all product lotteries  $\Delta(\mathcal{I}) \times \Delta(\mathcal{A})$ . We assume throughout that  $\succsim$  is complete, transitive, continuous and that its asymmetric part  $\succ$  is nonempty, so it admits a (nontrivial) continuous representation  $V : \Delta(\mathcal{I}) \times \Delta(\mathcal{A}) \rightarrow \mathbb{R}$ . That is, for any pair of product lotteries  $(\alpha, \ell)$  and  $(\alpha', \ell')$ ,  $(\alpha, \ell) \succsim (\alpha', \ell')$  if and only if  $V(\alpha, \ell) \geq V(\alpha', \ell')$ .

A society  $I$  is a finite subset of  $\mathcal{I}$ . Denote by  $\Delta(I)$  the set of identity lotteries over  $I$ .

**Definition 1.** *Utilitarianism:* The observer is a *utilitarian* if, for every society  $I \subset \mathcal{I}$ , her preferences restricted to  $\Delta(I) \times \Delta(\mathcal{A})$  admit a representation of the form

$$V(\alpha, \ell) = \sum_{i \in I} \alpha_i U_i(\ell_i)$$

where  $\ell_i \in \Delta(\mathcal{X})$  is the lottery faced by individual  $i$  (in which outcome  $a(i)$  is assigned a probability of  $\ell(a)$ ) and, for each individual  $i$ ,  $U_i(\ell_i) := \sum_{x \in \mathcal{X}} u_i(x) \ell_i(x)$  is an expected utility (EU) representation of  $\succsim_i$ .

As is well-known, the main behavioral property that characterizes EU preferences is *independence*:

**Definition 2.** *Independence:* Let  $\succsim^*$  be a preference relation on  $\Delta(\mathcal{X})$ . Then, for all  $p, q, r \in \Delta(\mathcal{X})$  and for all  $\beta \in [0, 1]$ ,

$$p \succ^* q \Rightarrow \beta p + (1 - \beta) r \succ^* \beta q + (1 - \beta) r$$

---

<sup>5</sup>We refrain from using the adjective ‘impartial’ until we make explicit, see Axiom 1 below, what we mean by it.

### 3 Assumptions

We make the following assumptions on  $\succsim$ :

**Axiom 1.** *Impartiality:* For any two individuals  $i, j \in \mathcal{I}$ ,

- (1) for all  $\ell \in \Delta(\mathcal{A})$ ,  $\succsim_i = \succsim_j$  and  $\ell_i = \ell_j \Rightarrow (i, \ell) \sim (j, \ell)$
- (2)  $(i, a_{\max}) \sim (j, a_{\max})$  and  $(i, a_{\min}) \sim (j, a_{\min})$

Part (1) of this axiom states that, given an action lottery  $\ell$ , if two individuals  $i$  and  $j$  with identical preferences are faced with the same action lottery, then the observer is indifferent between facing  $\ell$ , while being individual  $i$ , and facing  $\ell$ , while being individual  $j$ . This requirement seems quite natural. Part (2) says that being individual  $i$  and getting the most preferred outcome  $x_{\max}$  is assumed ethically equivalent to being individual  $j$  and getting the (same) most preferred outcome  $x_{\max}$ . As was convincingly explained by Karni (1998) who, in a different framework, employed a stronger axiom to derive utilitarianism “This value judgment ... is obtained by default. The methodological framework of revealed preference provides no ground for preferring one individual’s most preferred alternative over that of the other. Consequently, strict preference in either direction is either biased or involves considerations other than the rank order of the alternatives”. Clearly, the same applies to the worst outcome  $x_{\min}$ . A similar, though weaker, notion lies behind Segal’s (2000) *dictatorship indifference* axiom.<sup>6</sup>

Henceforth we assume that the observer preferences satisfy the impartiality axiom. To emphasize it, we call her an *impartial* observer.

**Axiom 2.** *Strong acceptance:* For all  $i \in \mathcal{I}$  and  $\ell, \ell' \in \Delta(\mathcal{A})$  satisfying  $\forall j \neq i \ell_j = \ell'_j$ , if  $\alpha_i > 0$  then

$$\ell_i \succsim_i \ell'_i \Leftrightarrow (\alpha, \ell) \succsim (\alpha, \ell')$$

This axiom states that the impartial observer sympathizes with individual  $i$  and fully adopts his preferences when she imagines herself being this individual with a positive probability, and when all other individuals are unaffected by her choice. This axiom strengthens

---

<sup>6</sup>See Karni (2003) for a different notion of impartiality.

Harsanyi's *acceptance* principle, according to which this sympathy to hold only when  $\alpha_i = 1$ . Axiom 2 also is equivalent to an axiom called *strong Pareto*, a version of Harsanyi's *Pareto* principle that was used in his aggregation analysis (see Harsanyi (1955), Weymark (1991) and Epstein and Segal (1992)).<sup>7</sup> In a sense, strong acceptance unifies two of Harsanyi's main ideas from his two famous analyses of social choice theory. The axiom is analogous to Karni's (1998) *sympathy* assumption.

The strong acceptance axiom enables us to express the impartial observer's function  $V$  as a *social welfare function*. That is,  $V$  can be expressed as a function  $W$  that, instead of the action lottery  $\ell$ , depends on the individuals' utilities associated with their induced lotteries  $\ell_i$ . More formally, let  $V_i(\ell_i) := V(i, \ell)$  be a representing utility the impartial observer attaches to individual  $i$  preferences. Note that, by impartiality,  $V_i(x_{\min}) = V_j(x_{\min}) := v_{\min}$  and  $V_i(x_{\max}) = V_j(x_{\max}) := v_{\max}$ , for all  $i, j \in \mathcal{I}$ , and hence by continuity, the image of  $V_i$ , for all  $i$ , is equal to the closed interval  $[v_{\min}, v_{\max}]$ . Then, strong acceptance implies that  $V(\alpha, \ell)$  can be written as  $W(\vec{\alpha}, \vec{V}(\ell))$ , where  $W$  is defined over  $\Delta([v_{\min}, v_{\max}])$ , the set of lotteries over all attainable utility values in which, for all  $i \in \text{Supp}(\alpha)$ ,  $\vec{\alpha}_i = \alpha_i$  is the probability of  $(\vec{V}(\ell))_i = V_i(\ell_i)$ . To see how  $W$  is constructed assume, for expositional clarity, that  $\text{Supp}(\alpha) = \{1, \dots, n\}$ . Then, given  $V$ ,  $V_i$  and  $\vec{v} = (v_1, \dots, v_n) \in [v_{\min}, v_{\max}]^n$ , define  $W$  by  $W(\vec{\alpha}, \vec{v}) := V(\alpha, \ell)$ , for the imaginary lottery  $\alpha$  satisfying  $\alpha_i = \vec{\alpha}_i$  and for any  $\ell$  satisfying  $v_i = V_i(\ell_i)$ , for all  $i \in \{1, \dots, n\}$ . By strong acceptance,  $W$  is well defined. Furthermore for a given  $\vec{\alpha}$ ,  $W$  is monotonic increasing with respect to  $v_i$  whenever  $\vec{\alpha}_i > 0$ . In the sequel we assume that  $W$  is normalized to satisfy  $W(1, v) = v$ , for all  $v \in [v_{\min}, v_{\max}]$ .

**Axiom 3.** *Convexity:* Consider two lotteries  $\ell, \ell' \in \Delta(\mathcal{A})$  for which there exist two individuals  $i$  and  $j$  satisfying  $\ell_i \succ_i \ell'_i$  and  $\ell_j \prec_j \ell'_j$ . Then, for any  $\alpha \in \Delta(\mathcal{I})$  satisfying  $\alpha_i > 0$ ,  $\alpha_j > 0$  and for any  $\beta \in (0, 1)$ ,

$$(\alpha, \ell) \succ (\alpha, \ell') \Rightarrow (\alpha, \beta\ell + (1 - \beta)\ell') \succ (\alpha, \ell')$$

That is, if two (participating) individuals disagree on the ranking of two action lotteries, then mixing the two lotteries cannot worsen the situation for the impartial observer. As was

---

<sup>7</sup>*Strong Pareto:* For a given society  $I$ , (1) for all lotteries  $\ell, \ell' \in \Delta(\mathcal{A})$ , if  $\ell_i \succ_i \ell'_i$  for all  $i$ , then  $\ell \succ \ell'$  and (2) if, furthermore, there exists an individual  $i'$  such that  $\ell_{i'} \succ_{i'} \ell'_{i'}$ , then  $\ell \succ \ell'$ .

explained in the introduction, this axiom is an expression of procedural fairness and is in accordance with Diamond's critique.

We include the requirement regarding two individuals with opposing preferences since procedural fairness has greater appeal when real conflict exists. Thus, our axiom is silent about situations where only one individual faces distinct lotteries under the action lotteries  $\ell$  and  $\ell'$ . As such, our axiom does not seem to imply that individual preferences must be convex but, as the reader may have noticed, this implication does hold when continuity is assumed.

Convexity is also related to social stability. Consider a society  $I \subset \mathcal{I}$ , whose set of available actions is given by a finite  $A \subset \mathcal{A}$ . For a given identity lottery  $\alpha \in \Delta(I)$ , the impartial observer's aim is to find the optimal action lottery that maximizes her utility. That is, the impartial observer seeks to solve the problem

$$\max_{\ell \in \Delta(A)} V(\alpha, \ell)$$

For societal stability, it is desirable that the set of optimal action lotteries does not change drastically when only minor changes occur. That is, we want this set to be upper hemi-continuous and convex valued with respect to the set of available actions  $A$ . Clearly, the continuity of  $\succsim$  implies upper hemi-continuity, while convexity is equivalent to the optimal set being a convex valued correspondence.

**Axiom 4.** *Indifference between identity and outcome lotteries:* For all societies  $\{i_1, \dots, i_n\}$  and for all sets of available actions  $\{a^1, \dots, a^n\}$ , if there exists  $k \in \{1, \dots, n\}$  such that  $(i_j, a^k) \sim (i_k, a^j)$  for all  $j$ , then

$$(\alpha^e, a^k) \sim (i_k, \ell^e)$$

where  $\alpha^e = \sum_{i=1}^n \frac{1}{n} i$  and  $\ell^e = \sum_{i=1}^n \frac{1}{n} a^i$ .

In words, this axiom states the following. Suppose there exists  $k$  such that the impartial observer is indifferent between facing the deterministic action  $a^k$ , imagining being individual  $i_j$ , and facing the deterministic action  $a^j$ , imagining being individual  $i_k$ , for all  $j$ . There are two ways to randomize over these degenerate pairs of equivalent product lotteries. The product lottery  $(\alpha^e, a^k)$  randomizes over identity lotteries (for the given action  $a^k$ ), while

product lottery  $(i_k, \ell^e)$  randomizes over action lotteries (for the given individual  $i_k$ ). Then, as implicitly assumed by Harsanyi, the impartial observer is indifferent between the two randomizations. Note that by continuity, the axiom holds for all  $\alpha$  (and for its corresponding action lottery  $\ell^\alpha$ ).

## 4 Utilitarianism

Our main result shows that the preceding axioms force all individuals to be of the EU type and, in addition, the impartial observer must be a utilitarian. That is, the behavioral assumptions on the impartial observer preferences induce her, as well as all individuals, to satisfy the independence axiom. This is achieved without imposing independence explicitly (neither on individuals nor on the observer).

**Theorem.** The impartial observer preferences satisfy *strong acceptance*, *convexity* and *indifference between identity and outcome lotteries* if, and only if, all individuals in  $\mathcal{I}$  satisfy *independence* and the impartial observer is a *utilitarian*.

The proof, which is relegated to the appendix, consists of two parts. First, we prove that all individuals in  $\mathcal{I}$  must satisfy the independence axiom. Then, we demonstrate that the impartial observer's preferences can be represented by a function that is additive with respect to the identity lotteries and that she, too, must satisfy the independence axiom.

**Comment 1.** Consider the Diamond example, described here by the table

	$a^1$	$a^2$
1	1	0
2	0	1

where individuals 1 and 2 correspond to the rows, actions  $a^1$  and  $a^2$  correspond to the columns and the entries represent the impartial observer's utilities. Considering the identity lottery  $\alpha^e = (\frac{1}{2}, \frac{1}{2})$ , action  $a^i$  corresponds to the pair  $(\alpha^e, a^i)$ , while tossing a fair coin corresponds to the pair  $(\alpha^e, \ell^e) = (\alpha^e, \frac{1}{2}a^1 + \frac{1}{2}a^2)$ . By indifference between identity and action lotteries,  $(\alpha^e, a^1) \sim (1, \ell^e)$  and  $(\alpha^e, a^2) \sim (2, \ell^e)$ . Hence,  $(1, \ell^e) \sim (2, \ell^e)$  and therefore, by strong

acceptance,  $(1, \ell^e) \sim (\alpha^e, \ell^e)$ . But then, by transitivity,  $(\alpha^e, a^1) \sim (\alpha^e, \ell^e)$  and the impartial observer is indifferent between the first action and the mixture. Put differently, she does not strictly prefer tossing a fair coin over the pure action  $a^1$ . Moreover, it can now be seen (proof omitted) that, by convexity, *any* mixture of the two actions  $a^1$  and  $a^2$  must be indifferent to  $a^1$ . This may seem like a significant step towards proving utilitarianism. However, the derivation of these ‘straight indifference line segments’ from the above extremely symmetric situation does not extend to the general case and cannot be utilized to derive a utilitarian representation.

**Comment 2.** As noted in the introduction, Blackorby, Donaldson and Mongin (2004) and Mongin and Pivato (2015) also derived utilitarianism without imposing independence. Although these authors work within Harsanyi’s aggregation theorem framework, a comparison to our theorem seems natural and is therefore, carried by focusing on the analysis of Mongin and Pivato (2015). Consider a given society  $I$ , with a set of actions  $A$ , and identify every product lottery  $(\alpha, \ell)$  with a matrix whose rows correspond to individuals and columns, to actions. Mongin and Pivato’s *ex ante* analysis is manifested by their *row preference* assumption, an assumption that is analogous to our strong acceptance axiom. Similarly, their *ex post* analysis is manifested by a *column preference* assumption that, in our model, would require an improvement in the impartial observer situation whenever an action  $a$  is replaced by a better action  $\bar{a}$ . Together with a *coordinate monotonicity* assumption, these two assumptions enable Mongin and Pivato to employ Gorman’s (1968) separability theorem and to derive a fully separable representation of the observer preferences. As can be seen in the appendix, our proof uses different arguments. Nevertheless, one might conjecture that indifference between identity and action lotteries must imply similar treatment of columns and rows and hence, together with strong acceptance, this means that Gorman’s separability theorem would yield our result. This, however, is not true. As can be seen in Examples 1 and 2 below, strong acceptance and indifference between identity and action lotteries are not sufficient to imply utilitarianism.

**Comment 3.** Another result that is close to ours appears in GKPS (2010). Their Theorem 3 roughly states that an impartial observer is a utilitarian if and only if she satisfies acceptance,

independence over identity lotteries and (their notion) of indifference between identity and action lotteries. However, as we do not assume any form of independence, the current result is stronger than theirs.<sup>8</sup>

The following first two examples demonstrate the necessity of convexity, while the third demonstrates the necessity of indifference between identity and action lotteries.

**Example 1.** Here we present a non utilitarian impartial observer who satisfies all axioms except for convexity. Assume that all preferences  $\succsim_i$  of individuals  $i \in \mathcal{I}$  belong to the rank-dependent utility class (RDU; see Weymark (1981) and Quiggin (1982)). Let  $g : [0, 1] \rightarrow [0, 1]$  be an increasing and onto function. For a given simple lottery  $r$  and  $z \in \text{Supp}(r)$  define  $F_r(z) := \sum_{y \leq z} r(y)$ ,  $F_r(z_-) := \sum_{y < z} r(y)$  and  $\nabla g(z; r) := g(F_r(z)) - g(F_r(z_-))$ . For simple lotteries, RDU preferences are represented by a function of the form  $V(p) = \sum_x u(x) \nabla g(x; p)$ . When  $g$  is the identity function,  $\nabla g(x; p) = p(x)$  and hence, in this case, RDU preferences are reduced to EU preferences. We assume that, in the eyes of the impartial observer, individual  $i$ 's preferences are represented by  $V_i(p) = \sum_x u_i(x) \nabla g(x; p)$ , where  $g$  is common to all individuals. The impartial observer preferences are also of the RDU type and are represented by

$$V^r(\alpha, \ell) = \sum_{i \in I} V_i(\ell_i) \nabla g(V_i(\ell_i); \alpha)$$

where the function  $g$  is the same function that used in the individuals' functions.

Impartiality and strong acceptance are satisfied by construction. To verify that indifference between identity and action lotteries is satisfied, consider, without loss of generality, a society  $I = \{1, \dots, n\}$ , a set of available actions  $\{a^1, \dots, a^n\}$  and assume there exists  $k$  for which  $V^r(j, a^k) = V^r(k, a^j)$  for all  $j$ . Then, for all  $j$ ,

$$u_j(a^k(j)) = V_j(a^k(j)) = V^r(j, a^k) = V^r(k, a^j) = V_k(a^j(k)) = u_k(a^j(k))$$

---

<sup>8</sup>It should also be noted that the notion of indifference between identity and action lotteries used by GKPS (2010) is stronger than ours. A formal claim (Claim 2), and its proof, appears in the appendix.

Hence,

$$\begin{aligned}
V^r(\alpha^e, a^k) &= \sum_{j \in I} u_j(a^k(j)) \nabla g(u_j(a^k(j)); \alpha^e) \\
&= \sum_{j \in I} u_k(a^j(k)) \nabla g(u_j(a^k(j)); \alpha^e) \\
&= \sum_{j \in I} u_k(a^j(k)) \nabla g(a^j(k); \ell_k^e) = V^r(k, \ell^e)
\end{aligned}$$

as required.

To see that convexity may fail to be satisfied, assume that  $g$  is strictly concave and consider  $j \in I$  and two distinct action lotteries  $\ell, \ell' \in \Delta(\mathcal{A})$  satisfying  $\ell_i = \ell'_i \forall i$ ,  $\ell_j \neq \ell'_j$  and  $V_j(\ell_j) = V_j(\ell'_j)$  (clearly, such lotteries exist). Then the strict concavity of  $g$  immediately implies  $V_j(\frac{1}{2}\ell_j + \frac{1}{2}\ell'_j) < V_j(\ell_j)$  and hence, for any  $\alpha$  with  $\alpha_j > 0$ ,  $V^r(\alpha, \frac{1}{2}\ell + \frac{1}{2}\ell') < V^r(\alpha, \ell)$ .

Note that, as the following case shows, non-convexity of  $\succsim_i$  (which is manifested by the concavity of  $g$ ), is not necessary for the non-convexity of  $\succsim$ . For this, let  $I = \{1, \dots, 5\}$ , consider the two actions described by the following matrix (the entries are the utility values)

	$a^1$	$a^2$
1	1	0
2	0	1
3	1	1
4	1	1
5	1	1

let  $g$  be given by the convex piecewise linear function

$$g(t) = \begin{cases} 0 & t \leq 0.2 \\ -\frac{1}{4} + \frac{5}{4}t & \text{otherwise} \end{cases}$$

and note that, by the convexity of  $g$ , each  $\succsim_i$  is convex.

Clearly,

$$V^r(\alpha^e, a^j) = g(0.2) \times 0 + (1 - g(0.2)) \times 1 = 1$$

for both  $j = 1, 2$ . Next, consider the lottery  $\frac{1}{2}a^1 + \frac{1}{2}a^2$ . For individuals 1 and 2,

$$V_i\left(\frac{1}{2}a^1(i) + \frac{1}{2}a^2(i)\right) = g(0.5) \times 0 + (1 - g(0.5)) \times 1 = \frac{5}{8}$$

hence, for the impartial observer,

$$\begin{aligned} V^r \left( \alpha^e, \frac{1}{2}a^1 + \frac{1}{2}a^2 \right) &= g(0.4) \times \frac{5}{8} + (1 - g(0.4)) \times 1 \\ &= \frac{1}{4} \times \frac{5}{8} + \frac{3}{4} \times 1 = \frac{29}{32} < 1 \end{aligned}$$

and convexity is not satisfied.

**Example 2.** In the two cases described in Example 1, either individual preferences are non-convex with respect to outcome lotteries (when  $g$  is concave) or the impartial observer preferences are non-convex with respect to identity lotteries (when  $g$  is convex). This might suggest that convexity would be satisfied if all preferences involved were convex. As we now show, this conjecture is false.

Assume that individual preferences are weighted utility (WU; see Chew (1983)). That is, for all  $i$  and  $p \in \Delta(\mathcal{X})$ ,

$$V_i(p) = V(p) = \sum_k p_k \frac{w(x_k)}{\sum_j p_j w(x_j)} u(x_k)$$

where  $u$  is a strictly increasing utility function and  $w$  is a non constant and positive weighting function. These preferences belong to the betweenness class (see Chew (1989) and Dekel (1986)), a class that is characterized by the property: for all lotteries  $p$  and  $q$ ,  $p \succcurlyeq q$  if and only if  $p \succcurlyeq \lambda p + (1 - \lambda)q \succcurlyeq q$ , for all  $\lambda \in (0, 1)$ . Clearly, betweenness implies that WU preferences are convex. Note that although individuals have identical preferences over  $\Delta(\mathcal{X})$ , they still may disagree on  $\Delta(\mathcal{A})$ .

The impartial observer preferences are of the same type and are given by

$$V^w(\alpha, \ell) = \sum_i \alpha_i \frac{w(u^{-1}(V(\ell_i)))}{\sum_j \alpha_j w(u^{-1}(V(\ell_j)))} V(\ell_i)$$

As in Example 1, indifference between identity and action lotteries is satisfied. To see it, assume (for  $k = 1$ )  $V^w(j, a^1(j)) = V^w(1, a^j(1))$ , for all  $j$ . That is,  $u(a^1(j)) = u(a^j(1))$  or, equivalently,  $a^1(j) = a^j(1)$ , for all  $j$ .

Then

$$\begin{aligned}
V^w(\alpha^e, a^1) &= \sum_i \frac{1}{n} \frac{w((u^{-1} \circ u)(a^1(i)))}{\sum_j \frac{1}{n} w((u^{-1} \circ u)(a^1(j)))} u(a^1(i)) \\
&= \sum_i \frac{1}{n} \frac{w(a^1(i))}{\sum_j \frac{1}{n} w(a^1(j))} u(a^1(i)) \\
&= \sum_i \frac{1}{n} \frac{w(a^i(1))}{\sum_j \frac{1}{n} w(a^j(1))} u(a^i(1)) \\
&= V(1, \ell^e)
\end{aligned}$$

Next we show that convexity is not satisfied. Consider again the Diamond's example and assume that  $w(u^{-1}(0)) = 0.1$ ,  $w(u^{-1}(\frac{2}{3})) = 0.6$ ,  $w(u^{-1}(\frac{32}{33})) = 0.75$  and  $w(u^{-1}(1)) = 0.8$ . Then,

$$V^w(\alpha^e, a^1) = \frac{\frac{1}{2}w(u^{-1}(1))}{\frac{1}{2}w(u^{-1}(1)) + \frac{1}{2}w(u^{-1}(0))} = \frac{0.5 \times 0.8}{0.5 \times 0.8 + 0.5 \times 0.1} = \frac{8}{9}$$

and

$$V^w(\alpha^e, a^2) = \frac{\frac{1}{2}w(u^{-1}(1))}{\frac{1}{2}w(u^{-1}(0)) + \frac{1}{2}w(u^{-1}(1))} = \frac{0.5 \times 0.8}{0.5 \times 0.1 + 0.5 \times 0.8} = \frac{8}{9}$$

Let  $\ell^p = 0.8a^1 + 0.2a^2$  be a mixture of  $a^1$  and  $a^2$ . Then,

$$\begin{aligned}
V(\ell_1^p) &= \frac{0.8w(u^{-1}(1))}{0.8w(u^{-1}(1)) + 0.2w(u^{-1}(0))} = \frac{0.8 \times 0.8}{0.8 \times 0.8 + 0.2 \times 0.1} = \frac{32}{33} \\
V(\ell_2^p) &= \frac{0.2w(u^{-1}(1))}{0.8w(u^{-1}(0)) + 0.2w(u^{-1}(1))} = \frac{0.2 \times 0.8}{0.8 \times 0.1 + 0.2 \times 0.8} = \frac{2}{3}
\end{aligned}$$

and, for the impartial observer,

$$\begin{aligned}
V^w(\alpha^e, \ell^p) &= \frac{1}{2} \frac{w(u^{-1}(V(\ell_1^p)))}{\frac{1}{2}w(u^{-1}(V(\ell_1^p))) + \frac{1}{2}w(u^{-1}(V(\ell_2^p)))} V(\ell_1^p) \\
&\quad + \frac{1}{2} \frac{w(u^{-1}(V(\ell_2^p)))}{\frac{1}{2}w(u^{-1}(V(\ell_1^p))) + \frac{1}{2}w(u^{-1}(V(\ell_2^p)))} V(\ell_2^p) \\
&= \frac{0.75}{0.75 + 0.6} \times \frac{32}{33} + \frac{0.6}{0.75 + 0.6} \times \frac{2}{3} \\
&\approx 0.835 < \frac{8}{9}
\end{aligned}$$

Hence, convexity is violated.

**Example 3.** A non utilitarian impartial observer who satisfies all axioms except for indifference between identity and action lotteries is the *generalized utilitarian* impartial observer of GKPS (2010). Consider

$$V^g(\alpha, \ell) = \sum_{i \in I} \alpha_i \phi_i [U_i(\ell_i)]$$

where  $\phi_i : [v_{\min}, v_{\max}] \rightarrow \mathbb{R}$  is strictly concave, for all  $i$ . It is easy to verify that strong acceptance and convexity are satisfied while, as was shown in GKPS, this observer deems identity lotteries inferior to action lotteries.

**Comment 4.** Consider the following assumption, which is weaker than the axiom on indifference between identity and action lotteries.

*Preference for identity lotteries:* For all societies  $\{i_1, \dots, i_n\}$  and for all sets of available actions  $\{a^1, \dots, a^n\}$ , if there exists  $k \in \{1, \dots, n\}$  such that  $(i_j, a^k) \sim (i_k, a^j)$  for all  $j$ , then

$$(\alpha^e, a^k) \succ (i_k, \ell^e)$$

In the appendix (Claim 1) we show that this assumption, in conjunction with convexity, implies indifference between identity and action lotteries. Therefore, our theorem could be stated in a stronger form. We prefer its current form because of the much greater normative appeal of the assumption of indifference between identity and action lotteries.

## 5 Conclusion

As stated in the introduction, we argue that, when societal decisions are at stake, postulates must be drawn from society-centered behavior. We have chosen to focus on the notion of procedural fairness (exhibited by convexity) and added to it the requirement that the impartial observer is indifferent between identity and action lotteries. In our main result we have shown that these two assumptions (together with strong acceptance) were sufficient to force the impartial observer to be a utilitarian. Unlike most utilitarian results, no form of the independence axiom was required here.

In addition to offering a society-centered basis for utilitarianism, our result sheds more light on what is needed in order to always have a strict preference for procedural fairness. According to our analysis, the latter implies that the impartial observer must display a preference for action lotteries. Two such non-utilitarian models exist in the literature. The first follows from Karni and Safra (2002).<sup>9</sup> In their model, which leads to the representation  $V(\alpha, \ell) = \sum_{i \in I} \alpha_i V_i(\ell_i)$ , individuals possess a sense of justice and the preference for procedural fairness is solely manifested by their behavior (their utilities  $V_i$  are assumed to be concave). It can easily be verified that this impartial observer displays a preference for action lotteries. The second model is the generalized utilitarian impartial observer of GKPS (2010). As mentioned above, GKPS show that a preference for action lotteries holds if and only if each  $\phi_i$  is concave, a condition that implies procedural fairness. For a third model, consider a rank dependent, or a Gini, impartial observer, whose preferences are represented by

$$V^{rd}(\alpha, \ell) = \sum_{i \in I} \phi(U_i(\ell_i)) \nabla g(U_i(\ell_i); \alpha)$$

(where each  $U_i$  is of the EU type and both  $\phi$  and  $g$  are concave). As can easily be verified, a preference for action lotteries follows from Chew, Karni and Safra (1987) while procedural fairness follows from Quiggin (1993, Section 9.1).

## 6 Appendix

**(A) Proof of the Theorem.** The ‘if’ part is immediate. The proof of the converse is divided into two parts.

(i)<sup>10</sup> In this part we show that all individuals satisfy the independence axiom. Consider an individual  $i^* \in \mathcal{I}$  and denote his preferences by  $\succsim^*$ . We want to demonstrate that for all  $p, q, r \in \Delta(\mathcal{X})$ ,  $p \sim^* q \Rightarrow \frac{1}{2}p + \frac{1}{2}r \sim^* \frac{1}{2}q + \frac{1}{2}r$ . This, using Herstein and Milnor (1953), is suffice to prove that  $\succsim^*$  satisfies the independence axiom. Using the continuity of  $\succsim^*$ , we can restrict attention to equi-probability lotteries with the same number of outcomes:

<sup>9</sup>See also Grant, Kajii, Polak and Safra (2012).

<sup>10</sup>The proof of this part is similar to that of Dekel, Safra and Segal (1991, Theorem 2). However, dealing with social multi-person frameowrk, our proof is more general than (and improves upon) theirs.

$p = ((\frac{1}{k}, \dots, \frac{1}{k}), x)$ ,  $q = ((\frac{1}{k}, \dots, \frac{1}{k}), y)$ , and  $r = ((\frac{1}{k}, \dots, \frac{1}{k}), z)$  where, in addition and without loss of generality,  $x, y$  and  $z$  are non-constant vectors and  $p \succ^* r$ . Without loss of generality, assume that  $x_1 > x_j$  and  $y_1 > y_j$ , for all  $j \neq 1$ .

Consider a society  $I$  consisting of  $n = 2k$  individuals, all with preferences  $\succsim_i = \succsim^*$  and assume, without loss of generality, that  $I = \{1, \dots, n\}$ .<sup>11</sup> Denote by  $\Pi^n$  the set of  $n$  permutations on  $\{1, \dots, n\}$  of the form  $\pi_1 = (1, 2, \dots, n)$ ,  $\pi_2 = (2, 3, \dots, 1)$ , ...,  $\pi_n = (n, 1, 2, \dots, n-1)$  (where  $\pi_j(i)$  stands for the  $i$ th element of the permutation  $\pi_j$ ). We concentrate on a set of actions  $\dot{A} = \{\dot{a}^1, \dots, \dot{a}^n\}$  available to the society that are defined as follows: for  $j = 1, \dots, k$

$$\dot{a}^j(i) = \begin{cases} x_{\pi_j(i)} & \text{if } 1 \leq i \leq k \\ z_{\pi_j(i-k)} & \text{if } k < i \leq n \end{cases}$$

and, for  $j = k+1, \dots, n$

$$\dot{a}^j(i) = \begin{cases} z_{\pi_{j-k}(i)} & \text{if } 1 \leq i \leq k \\ x_{\pi_{j-k}(i-k)} & \text{if } k < i \leq n \end{cases}$$

(i.e.,  $\dot{a}^1 = (x_1, x_2, \dots, x_k, z_1, z_2, \dots, z_k)$ ,  $\dot{a}^2 = (x_2, x_3, \dots, x_k, x_1, z_2, z_3, \dots, z_k, z_1)$ , . . . ,  $\dot{a}^{k+1} = (z_1, z_2, \dots, z_k, x_1, x_2, \dots, x_k)$ ,  $\dot{a}^{k+2} = (z_2, z_3, \dots, z_k, z_1, x_2, x_3, \dots, x_k, x_1)$ , etc.).

We start by showing that, for all  $i$ ,  $(\alpha^e, \dot{a}^i) \sim (\alpha^e, \ell^e)$ . For  $i = 1$ , since in both  $(j, \dot{a}^1)$  and  $(1, \dot{a}^j)$  the impartial observer faces the same deterministic outcome ( $x_j$  if  $j \leq k$  and  $z_{j-k}$  otherwise) then, by impartiality,  $(j, \dot{a}^1) \sim (1, \dot{a}^j)$ , for all  $j \in I$ . By indifference between identity and action lotteries,  $(\alpha^e, \dot{a}^1) \sim (1, \ell^e)$ . Since for all  $i, j$ ,  $\ell_i^e = \ell_j^e$ , impartiality implies  $(1, \ell^e) \sim (\alpha^e, \ell^e)$  and hence, by transitivity,  $(\alpha^e, \dot{a}^1) \sim (\alpha^e, \ell^e)$ . By similar arguments,  $(\alpha^e, \dot{a}^i) \sim (\alpha^e, \ell^e)$  for all  $i$ .

Next, we show that the action lottery  $\ell^k = \frac{1}{k} \sum_{j=1}^k \dot{a}^j$  satisfies  $(\alpha^e, \ell^k) \sim (\alpha^e, \ell^e)$ . First note that, by construction, individual  $i$  strictly prefers action  $\dot{a}^i$  over all other actions and, by monotonicity with respect to first-order stochastic-dominance, he also strictly prefers action  $\dot{a}^i$  over all mixtures of the other actions. Hence, by repeated application of convexity,  $(\alpha^e, \ell^k) \succ (\alpha^e, \dot{a}^j)$  for all  $j = 1, \dots, k$  and, therefore,  $(\alpha^e, \ell^k) \succ (\alpha^e, \ell^e)$ . For the converse, consider the action lottery  $\hat{\ell}^k = \frac{1}{k} \sum_{j=k+1}^n \dot{a}^j$  and note that, for all  $i = 1, \dots, k$ ,  $\hat{\ell}_i^k$ , the lottery

---

<sup>11</sup>Here we use the infinity of the set  $\mathcal{I}$ .

individual  $i$  faces under  $\hat{\ell}^k$ , is identical to  $\ell_{k+i}^k$ , the lottery that individual  $k+i$  faces under  $\ell^k$ , and  $\hat{\ell}_{k+i}^k$ , the lottery individual  $k+i$  faces under  $\hat{\ell}^k$ , is identical to  $\ell_i^k$ , the lottery that individual  $i$  faces under  $\ell^k$ . Hence, by impartiality,  $(\alpha^e, \hat{\ell}^k) \sim (\alpha^e, \ell^k)$ . By construction, individual 1 strictly prefers  $\ell^k$  (where he faces the lottery  $p$ ) over  $\hat{\ell}^k$  (where he faces the lottery  $r$ ), while individual  $k+1$  has the opposite preferences. Therefore, as  $\ell^e = \frac{1}{2}\hat{\ell}^k + \frac{1}{2}\ell^k$ , convexity implies  $(\alpha^e, \ell^e) \succ (\alpha^e, \ell^k)$ . Hence,  $(\alpha^e, \ell^k) \sim (\alpha^e, \ell^e)$ . By impartiality and transitivity, we then get  $(\alpha^e, \ell^k) \sim (1, \ell^e)$ . Note that in the first lottery, the first  $k$  individuals face the lottery  $p$  and the rest face the lottery  $r$  while, in the latter lottery, individual 1 is faced with the lottery  $\frac{1}{2}p + \frac{1}{2}r$ .

Finally, consider a second society with the same set of individuals  $I$  and a set of available actions  $\tilde{A} = \{\tilde{a}^1, \dots, \tilde{a}^{2k}\}$  that is derived from  $\dot{A}$  by replacing every  $x_j$  by  $y_j$ . Clearly, a similar conclusion holds: the impartial observer is indifferent between the product lottery  $(\alpha^e, \ell^k)$ , in which the first  $k$  individuals face the lottery  $q$  and the rest face the lottery  $r$ , and the product lottery  $(1, \ell^e)$ , in which individual 1 is faced with the lottery  $\frac{1}{2}q + \frac{1}{2}r$ . But as  $p \sim^* q$ , all individuals in  $I$  are indifferent between  $p$  and  $q$  and hence, by strong acceptance, the impartial observer is indifferent between the two occurrences (in the two societies) of the product lottery  $(\alpha^e, \ell^k)$ . Thus by transitivity, the impartial observer, while imagining herself being individual 1, is indifferent between facing the lotteries  $\frac{1}{2}p + \frac{1}{2}r$  and  $\frac{1}{2}q + \frac{1}{2}r$ . By strong acceptance,  $\frac{1}{2}p + \frac{1}{2}r \sim^* \frac{1}{2}q + \frac{1}{2}r$  and hence  $\succ^*$  satisfies independence.

(ii) In the second part we show that the impartial observer is a utilitarian. Consider a society  $I$  (without loss of generality,  $I = \{1, \dots, n\}$ ) and let  $V(\alpha, \ell) = W(\vec{\alpha}, \vec{V}(\ell))$  be a representation of the impartial observer preferences where  $(\vec{V}(\ell))_i = V_i(\ell_i) = \varphi_i(U_i(\ell_i))$ ,  $\varphi_i$  is monotonic increasing and, by part (ii),  $U_i(\ell_i) = \sum_{x \in \mathcal{X}} u_i(x) \ell_i(x)$  is an EU representation of individual  $i$ 's preferences. Since  $u_i$  is determined up to (positive) affine transformations, we can assume it satisfies  $u_i(x_{\min}) = v_{\min}$  and  $u_i(x_{\max}) = v_{\max}$  (hence,  $\varphi_i(v_{\min}) = v_{\min}$  and  $\varphi_i(v_{\max}) = v_{\max}$ , for all  $i$ ). Choose  $(\alpha, \ell) \in \Delta(I) \times \Delta(\mathcal{A})$ , denote  $v_i = \varphi_i(U_i(\ell_i))$  and let  $c_i(\ell_i) \in \mathcal{X}$  be individual  $i$ 's certainty equivalent of the lottery  $\ell_i$  (that is,  $u_i(c_i(\ell_i)) = U_i(\ell_i)$ ). Consider a set of actions  $\hat{A} = \{\hat{a}^j \mid j \in \{1, \dots, n\}\}$  satisfying  $\hat{a}^1(i) = c_i(\ell_i)$  and  $\hat{a}^j(1) = (\varphi_1 \circ u_1)^{-1}(v_j)$  for  $i, j = 1, \dots, n$ . By construction,

$V(i, \hat{a}^1) = (\varphi_i \circ u_i)(c_i(\ell_i)) = v_i$  and  $V(1, \hat{a}^i) = (\varphi_1 \circ u_1) \circ (\varphi_1 \circ u_1)^{-1}(v_i) = v_i$ . Hence  $(i, \hat{a}^1) \sim (1, \hat{a}^i)$  and, by indifference between identity and action lotteries,  $(\alpha, \hat{a}^1) \sim (1, \ell^\alpha)$  ( $\ell^\alpha$  is the action lottery associated with  $\alpha$ ). Put differently,  $V(\alpha, \hat{a}^1) = V(1, \ell^\alpha)$  or, equivalently,  $W(\vec{\alpha}, (v_1, \dots, v_n)) = W(1, \varphi_1(U_1(\ell_1^\alpha)))$ . Therefore,

$$\begin{aligned} V(\alpha, \ell) &= W(\vec{\alpha}, \vec{V}(\ell)) = W(\vec{\alpha}, (v_1, \dots, v_n)) = W(1, \varphi_1(U_1(\ell_1^\alpha))) \\ &= \varphi_1(U_1(\ell_1^\alpha)) = \varphi_1\left(\sum_{i=1}^n \alpha_i u_i((\varphi_1 \circ u_1)^{-1}(v_i))\right) \\ &= \varphi_1\left(\sum_{i=1}^n \alpha_i \varphi_1^{-1}(v_i)\right) = \varphi_1\left(\sum_{i=1}^n \alpha_i (\varphi_1^{-1} \circ \varphi_i)(U_i(\ell_i))\right) \end{aligned}$$

Denote  $\bar{V} = \varphi_1^{-1} \circ V$  and  $\phi_i = \varphi_1^{-1} \circ \varphi_i$  (note that  $\bar{V}$  also represents the impartial observer preferences and its image is  $[v_{\min}, v_{\max}]$ ). By the above,

$$\bar{V}(\alpha, \ell) = \sum_{i=1}^n \alpha_i \phi_i[U_i(\ell_i)]$$

To conclude, we show that for all  $i$ ,  $\bar{V}_i = \phi_i \circ U_i$  is affine which, given  $\varphi_i(v_{\min}) = v_{\min}$  and  $\varphi_i(v_{\max}) = v_{\max}$ , implies  $\bar{V}_i = U_i$ . Take  $\ell, \ell' \in \Delta(\mathcal{A})$ . Since  $\succsim_i$  is of the EU type, we have that for all  $\lambda \in [0, 1]$ ,

$$\begin{aligned} \bar{V}_i(\lambda \ell_i + (1 - \lambda) \ell'_i) &= \phi_i[U_i(\lambda \ell_i + (1 - \lambda) \ell'_i)] = \phi_i[\lambda U_i(\ell_i) + (1 - \lambda) U_i(\ell'_i)] \\ &= \phi_i[\lambda u_i(c_i(\ell_i)) + (1 - \lambda) u_i(c_i(\ell'_i))] = \phi_i[U_i(\lambda c_i(\ell_i) + (1 - \lambda) c_i(\ell'_i))] \\ &= \bar{V}_i(\lambda \check{a}^i(i) + (1 - \lambda) \check{a}^j(i)) = \bar{V}(i, \lambda \check{a}^i + (1 - \lambda) \check{a}^j) \end{aligned}$$

for actions  $\check{a}^i$  and  $\check{a}^j$  satisfying  $\check{a}^i(i) = c_i(\ell_i)$ ,  $\check{a}^j(i) = c_i(\ell'_i)$  and  $\check{a}^i(j) = (\varphi_j \circ u_j)^{-1} \circ (\varphi_i \circ u_i)(c_i(\ell'_i))$  (note that the element  $\lambda c_i(\ell_i) + (1 - \lambda) c_i(\ell'_i)$  that appears in the second line is a lottery, not an outcome). By construction,

$$\bar{V}(j, \check{a}^j) = (\varphi_j \circ u_j) \circ (\varphi_j \circ u_j)^{-1} \circ (\varphi_i \circ u_i)(c_i(\ell'_i)) = (\varphi_i \circ u_i)(c_i(\ell'_i)) = \bar{V}(i, \check{a}^j)$$

and hence, by indifference between identity and action lotteries and for  $\alpha$  satisfying  $\alpha_i = \lambda$ ,  $\alpha_j = 1 - \lambda$  and  $\alpha_k = 0$  otherwise,

$$\bar{V}(i, \lambda \check{a}^i + (1 - \lambda) \check{a}^j) = \bar{V}(\lambda i + (1 - \lambda) j, \check{a}^i)$$

(note that actions  $\check{a}^k$  for  $k \neq i, j$  are irrelevant but can easily be defined as to fit with the requirements of the axiom). Now, by the structure of  $\bar{V}$ ,

$$\begin{aligned}\bar{V}(\lambda i + (1 - \lambda)j, \check{a}^i) &= \lambda \bar{V}(i, \check{a}^i) + (1 - \lambda) \bar{V}(j, \check{a}^i) = \lambda \bar{V}(i, \check{a}^i) + (1 - \lambda) \bar{V}(i, \check{a}^j) \\ &= \lambda \bar{V}_i(\check{a}^i(i)) + (1 - \lambda) \bar{V}_i(\check{a}^j(i)) = \lambda \bar{V}_i(c_i(\ell_i)) + (1 - \lambda) \bar{V}_i(c_i(\ell'_i)) \\ &= \lambda \bar{V}_i(\ell_i) + (1 - \lambda) \bar{V}_i(\ell'_i)\end{aligned}$$

Summarizing,

$$\bar{V}_i(\lambda \ell_i + (1 - \lambda) \ell'_i) = \lambda \bar{V}_i(\ell_i) + (1 - \lambda) \bar{V}_i(\ell'_i)$$

and the affinity of  $\bar{V}_i$  is established. ■

### (B) Preference for identity lotteries vs indifference between identity and action lotteries.

**Claim 1.** If the impartial observer preferences satisfy *strong acceptance*, *convexity* and *preference for identity lotteries* than they satisfy *indifference between identity and outcome lotteries*.

**Proof.** Consider, without loss of generality, a society  $I = \{1, \dots, n\}$ , a set of available actions  $A = \{a^1, \dots, a^n\}$  and assume that (again, without loss of generality)  $V(i, a^1) = V(1, a^i) := v_i$ , for all  $i$ . Without loss of generality we can assume that all  $v_i$  are pairwise different and that  $v_i > v_{i+1}$  for all  $i < n$ . For  $i, j \in \{1, \dots, n\}$ , let  $x_{ij} \in \mathcal{X}$  be defined by  $V_i(x_{ij}) = v_{\pi_j(i)}$ , where  $\pi_j \in \Pi^k$ , and note that, by the monotonicity of each  $V_i$  with respect to the outcomes of  $\mathcal{X}$ ,  $V_1(x_{11}) > V_1(x_{12}) > \dots > V_1(x_{1n})$ ,  $V_2(x_{2n}) > V_2(x_{21}) > V_2(x_{22}) > \dots > V_2(x_{2(n-1)})$ , ...,  $V_n(x_{n2}) > V_n(x_{n3}) > \dots > V_n(x_{nn}) > V_n(x_{n1})$ . Consider a new set of actions  $\bar{A} = \{\bar{a}^1, \dots, \bar{a}^n\}$  satisfying  $\bar{a}^j(i) = x_{ij}$ . By construction,

$$V(i, \bar{a}^1) = V_i(x_{i1}) = v_{\pi_1(i)} = v_i = V(i, a^1)$$

and

$$V(1, \bar{a}^i) = V_1(x_{1i}) = v_{\pi_i(1)} = v_i = V(1, a^i)$$

which implies that, by strong acceptance,  $V(\alpha^e, a^1) = W(\alpha^e, (v_1, \dots, v_n)) = V(\alpha^e, \bar{a}^1)$  and  $V(1, \ell^e)$ , given  $A$ , is equal to  $V(1, \ell^e)$ , given  $\bar{A}$ . Hence it is sufficient to restrict attention to  $\bar{A}$  and to show that  $V(\alpha^e, \bar{a}^1) = V(1, \ell^e)$  (given  $\bar{A}$ ). For this note that: (i)

since  $V(\alpha^e, \bar{a}^i) = W(\alpha^e, (v_1, \dots, v_n))$  for all  $i$ , we have  $V(\alpha^e, \bar{a}^i) = V(\alpha^e, \bar{a}^j)$ , for all  $i, j$ ; (ii) by construction, for every  $k \in \{1, \dots, n\}$ ,  $V(i, \bar{a}^k) = V(k, \bar{a}^i)$ , for all  $i$ ; (iii)  $V(\alpha^e, \ell^e) \in [\min_i V(i, \ell^e), \max_i V(i, \ell^e)]$  and hence, if  $V(\alpha^e, \ell^e) = \max_i V(i, \ell^e)$  then  $V(\alpha^e, \ell^e) = V(i, \ell^e) = V(j, \ell^e)$ , for all  $i, j$ ; and (iv) individual  $i$  strictly prefers action  $\bar{a}^{n+2-i}$  (where  $\bar{a}^{n+2-1} = \bar{a}^{n+1} := \bar{a}^1$ ) over all other actions and, by the monotonicity of  $V_i$  with respect to first-order stochastic-dominance, he strictly prefers action  $\bar{a}^i$  over all mixtures of the other actions. Therefore,

$$V(\alpha^e, \bar{a}^1) = \max_k V(\alpha^e, \bar{a}^k) \geq \max_k V(k, \ell^e) \geq V(\alpha^e, \ell^e) \geq V(\alpha^e, \bar{a}^1)$$

where the equality follows from (i), the first inequality follows from (ii) and from preference for identity lotteries, the second inequality follows from the first part of (iii) and the last inequality follows from (iv) by repeated application of convexity (note that  $\ell^e = \frac{1}{n} \sum_j \bar{a}^j$ ).

Since the first and the last elements are identical,  $\max_k V(k, \ell^e) = V(\alpha^e, \ell^e)$  which, by the second part of (iii), implies that  $V(1, \ell^e) = \max_k V(k, \ell^e)$  and, therefore,  $V(1, \ell^e) = V(\alpha^e, \bar{a}^1)$ . Hence the impartial observer is indifferent between identity and action lotteries. ■

**(C) GKPS's (2010) indifference between identity and action lotteries implies ours.**

**Claim 2.** As in GKPS (2010), assume that  $\forall \alpha, \alpha' \in \Delta(\mathcal{I}), \forall \ell, \ell' \in \Delta(\mathcal{A})$  and  $\forall \beta \in (0, 1)$ ,

$$(\alpha, \ell') \sim (\alpha', \ell) \Rightarrow (\beta\alpha + (1 - \beta)\alpha', \ell) \sim (\alpha, \beta\ell + (1 - \beta)\ell')$$

Then the impartial observer exhibits indifference between identity and action lotteries.

**Proof.** The proof is by induction. Without loss of generality, consider a society  $I = \{1, \dots, n\}$ , the set of available actions  $A = \{a^1, \dots, a^n\}$  and assume that  $(1, a^i) \sim (i, a^1)$ , for all  $i$ .

First let  $n = 2$ . By the GKPS. condition,  $(1, a^2) \sim (2, a^1)$  implies

$$\left(\frac{1}{2}1 + \frac{1}{2}2, a^1\right) \sim \left(1, \frac{1}{2}a^1 + \frac{1}{2}a^2\right)$$

as required.

Next assume it holds for  $n - 1$  and consider  $n$ . Assume, without loss of generality, that the acts of  $A$  satisfy  $(i, a^j) \sim (i + 1, a^{j-1})$  for all  $i \in \{1, \dots, n - 1\}$ ,  $j \in \{2, \dots, n\}$ . Consider the society  $I^{\setminus 1} = \{2, \dots, n\}$  and the set of actions  $A^{\setminus n} = \{a^1, \dots, a^{n-1}\}$ . By construction,  $(2, a^i) \sim (i + 1, a^1)$  for all  $i = 1, \dots, n - 1$  and hence, by the induction hypothesis,  $(\frac{1}{n-1} \sum_{i=2}^n i, a^1) \sim (2, \frac{1}{n-1} \sum_{i=1}^{n-1} a^i)$ . Next apply the same argument to  $I^{\setminus n} = \{1, \dots, n - 1\}$  and  $A^{\setminus n} = \{a^1, \dots, a^{n-1}\}$ , where  $(2, a^i) \sim (i, a^2)$  for all  $i$ , to get  $(2, \frac{1}{n-1} \sum_{i=1}^{n-1} a^i) \sim (\frac{1}{n-1} \sum_{i=1}^{n-1} i, a^2)$ . Finally, apply it to  $I^{\setminus n} = \{1, \dots, n - 1\}$  and  $A^{\setminus 1} = \{a^2, \dots, a^n\}$ , where  $(1, a^{i+1}) \sim (i, a^2)$  for all  $i$ , to get  $(\frac{1}{n-1} \sum_{i=1}^{n-1} i, a^2) \sim (1, \frac{1}{n-1} \sum_{i=2}^n a^i)$ . By transitivity,

$$(\frac{1}{n-1} \sum_{i=2}^n i, a^1) \sim (1, \frac{1}{n-1} \sum_{i=2}^n a^i)$$

To conclude, mix both sides of the last indifference with  $(1, a^1)$  and, by the GKPS. condition, obtain  $(\alpha^e, a^1) \sim (1, \ell^e)$  for  $I = \{1, \dots, n\}$  and  $A = \{a^1, \dots, a^n\}$ , as required.  $\blacksquare$

## References

- [1] Blackorby, C., D. Donaldson and P. Mongin. 2004. Social aggregation without the expected utility hypothesis. Discussion paper 2004-020, Ecole Polytechnique.
- [2] Chew, S.H. 1983. A generalization of the quasilinear mean with applications to the measurement of income inequality and decision theory resolving the Allais paradox. *Econometrica*, 51, 1065-1092.
- [3] Chew, S.H. 1989. Axiomatic utility theories with the betweenness property. *Annals of operations Research*, 19, 273-298.
- [4] Chew, S.H., L. Epstein and U. Segal. 1991. Mixture symmetry and quadratic utility. *Econometrica*, 59, 139-163.
- [5] Chew, S.H., E. Karni and Z. Safra. 1987. Risk aversion in the theory of expected utility with rank dependent probabilities. *Journal of Economic Theory*, 42, 370-381.
- [6] Dekel, E. 1986. An axiomatic characterization of preferences under uncertainty: weakening the independence axiom. *Journal of Economic Theory*, 40, 304-318.

- [7] Dekel, E., Z. Safra and U. Segal. 1991. Existence and dynamic consistency of Nash equilibrium with non-expected utility preferences. *Journal Economic Theory*, 55, 229-246.
- [8] Dhillon, A. and J.F. Mertens. 1999. Relative utilitarianism. *Econometrica*, 67(3), 471-498.
- [9] Diamond, P. A. 1967. Cardinal welfare, individualistic ethics, and interpersonal comparison of utility: Comment. *The Journal of Political Economy*, 75(5), 765-766.
- [10] Elster, J. 1989. *Solomonic Judgements: Studies in the Limitation of Rationality*. Cambridge University Press.
- [11] Epstein, L. and U. Segal. 1992. Quadratic social welfare functions. *Journal of Political Economy*, 100(4), 691-712.
- [12] Fleurbaey, M. and P. Mongin. 2012. The utilitarian relevance of the aggregation theorem. Mimeo.
- [13] Gilboa, I., D. Samet and D. Schmeidler. 2004. Utilitarian aggregation of beliefs and tastes. *Journal of Political Economy*, 112(4), 932-938.
- [14] Gorman, W.M. 1968. The structure of utility functions. *Review of Economic Studies*, 35, 369-390.
- [15] Grant, S., A. Kajii, B. Polak and Z. Safra. 2010. General utilitarianism and Harsanyi's impartial observer theorem. *Econometrica*, 78(6), 1939-1971.
- [16] Grant, S., A. Kajii, B. Polak and Z. Safra. 2012. A generalized representation theorem for Harsanyi's ('impartial') observer. *Social Choice and Welfare*, 39, 833-846.
- [17] Harsanyi, J.C. 1953. Cardinal utility in welfare economics and in the theory of risk-taking. *Journal of Political Economy*, 61, 434-435.
- [18] Harsanyi, J.C. 1955. Cardinal welfare, individualistic ethics, and interpersonal comparisons of utility. *Journal of Political Economy*, 63, 309-321.

- [19] Harsanyi, J.C. 1975. Nonlinear social welfare functions. *Theory and Decision*, 6, 311-332.
- [20] Harsanyi, J.C. 1977. *Rational Behavior and Bargaining Equilibrium in Games and Social Situations*. Cambridge University Press.
- [21] Herstein, I. N. and J. Milnor. 1953. An axiomatic approach to measurable utility. *Econometrica*, 21(2), 291-297.
- [22] Karni, E. 1998. Impartiality: definition and representation. *Econometrica*, 66(6), 1405-1415.
- [23] Karni, E. 2003. Impartiality and interpersonal comparisons of variations in well-being. *Social Choice and Welfare*, 21, 95-111.
- [24] Karni, E and Z. Safra. 2002. Individual sense of justice: a utility representation. *Econometrica*, 70, 263-284.
- [25] Mongin, P. and M. Pivato. 2015. Ranking multidimensional alternatives and uncertain prospects. *Journal of Economic Theory*, 157, 146-171.
- [26] Quiggin J. 1982. A theory of anticipated utility. *Journal of Economic Behavior and Organization*, 3, 323-343.
- [27] Quiggin J. 1993. *Generalized expected utility theory: The rank-dependent model*. Kluwer Academic Publishers.
- [28] Segal, U. 2000. Let's agree that all dictatorships are equally bad. *Journal of Political Economy*, 108(3), 569-589.
- [29] Sen, A.K. 1976. Welfare inequalities and Rawlsian axiomatics. *Theory and Decision*, 7, 243-262.
- [30] Sen, A.K. 1977. Non-linear social welfare functions: A reply to Professor Harsanyi. In *Foundational Problems in the Social Sciences*, R. Butts and J. Hintikka (eds.). Reidel Publishing Company, 297-302.

- [31] Weymark, J.A. 1981. Generalized Gini inequality indices. *Mathematical Social Sciences*, 1(4), 409-430.
- [32] Weymark, J.A. 1991. A reconsideration of the Harsanyi-Sen debate on utilitarianism. In *Interpersonal Comparisons of Well-being*, J. Elster and J.E. Roemer (eds). Cambridge University Press, 255-320.
- [33] Zhou, L. 1997. Harsanyi's utilitarianism theorems: general societies. *Journal of Economic Theory*, 72(1), 198-207.