

Similarity and Moral Trade-offs

Alex Voorhoeve (corresponding author)

Department of Philosophy, Logic and Scientific Method

London School of Economics and Political Science

Houghton Street, London, WC2A 2AE, UK

a [dot] e [dot] Voorhoeve [at-sign] lse [dot] ac [dot] uk

&

Department of Bioethics

U.S. National Institutes of Health

Arnaldur Stefansson

Department of Economics

Uppsala University

Brian Wallace

Department of Economics

University College London

Author note: Brian Wallace is now at Source Clear, San Francisco.

Interests: None of the authors have any interests to declare.

Word count (incl. tables and references): 9169

Abstract

How, and how reliably, do people make difficult moral trade-offs? We pursue this question through an experiment in which subjects must either save a larger number of people from a smaller harm or save a smaller number of people from a greater harm. Our results indicate use of a similarity heuristic by around two-fifths of subjects. When alternatives appear dissimilar in terms of the number of people that can be saved but similar in terms of the magnitude of harm from which they can be saved, this heuristic mandates saving the greater number. In our experiment, use of this heuristic leads to violations of principles of rational choice at the individual and collective level. It also leads to choices that are inconsistent with all standard theories of distributive justice. We argue that this demonstrates the unreliability of moral judgments in cases that elicit similarity-based choice.

Keywords: Similarity, moral decision-making, heuristics, biases, moral psychology.

Introduction

How do people make difficult moral trade-offs, such as when they must decide whether to save twenty people from a moderate harm or instead five other people from a large harm? This question is, of course, relevant for the development of descriptive theories of choice. But it is also pertinent for public policy and moral philosophy. For example, debates about the use of public resources for health are often informed by surveys of the public's views on such trade-offs (see, e.g. Ubel et al. 1996; Gaertner and Schokkaert 2012, chapter 5; Nord and Johansen 2014). And in moral philosophy, people's opinions about moral cases are used in the search for plausible moral principles that explain and justify confidently held case judgments (Rawls 1999; Daniels 2013). For both policy and philosophy, it is therefore important to establish under which circumstances people's case judgments reflect the considered values of the people surveyed. Understanding the process by which people arrive at their decisions can help us answer this question by establishing whether individuals use particular heuristics and whether these heuristics make them smart (cf. Gigerenzer, Todd, et al. 2000) or instead generate objectionable biases (cf. Tversky and Kahneman 1974).

In this paper, we aim to contribute to these empirical and normative projects. We focus on cases in which individuals take on the role of a decision-maker about the use of public health care resources and are asked to balance the number of people they can save against the magnitude of the harm from which they can save them. Subjects presumably hold values that could inform a theory of justice for such cases, which is why social scientists conduct surveys and philosophers devise thought experiments to uncover them. However, subjects are unlikely to have at the ready a fully developed theory of justice to decide them. Moreover, they are unlikely to have much experience with trading off the magnitude of

harm prevented against the number of people saved from harm. They can therefore be expected to find some such trade-offs difficult to make.

Some descriptive theories of choice hold that when faced with challenging choices between two-dimensional alternatives, a substantial share of people first see if they can use a heuristic consisting of a simple rule (or a set of such rules, sequentially applied) to make a choice without explicitly trading off dimensions against each other (Tversky 1972, Brandstätter et al. 2006, Manzini and Mariotti 2007, Drechsler et al. 2014, Tserenjigmid 2015). In this paper, we investigate the use of one such heuristic, known as “similarity-based decision-making” (Tversky 1969, Rubinstein 1988). Subjects who use this heuristic decide as follows when faced with a pair of two-dimensional alternatives. If the alternatives are similar along one dimension and dissimilar along another, they choose the alternative that is better along the dissimilar dimension.

We report an experiment designed to test for use of this heuristic. Our results suggest that somewhat in excess of 40% of subjects employ it. Moreover, use of this heuristic induces individual and collective choices that are inconsistent with both formal theories of rationality and all standard, substantive theories of distributive justice. We argue that these results indicate the unreliability of moral judgments in cases that elicit similarity-based decision-making.

We proceed as follows. In section 1, we describe the similarity heuristic in more detail and review evidence of its use. In section 2, we describe the idea underlying our experiment. In section 3, we describe our methods. In section 4, we discuss our results. In section 5, we state our principal empirical and normative conclusions.

1. Similarity-based decision-making: theory and previous experiments

A general hypothesis about the role of similarity in a pairwise choice between multi-dimensional alternatives runs as follows. Dimensions along which alternatives appear similar will receive less attention and so receive less weight, while dimensions along which they appear dissimilar will capture attention and so receive greater weight (Mellers and Biagini 1994, Goldstone et al. 1997, Dhar et al. 1999, Köszegi and Szeidl 2013). Here, we focus on the following version of this general hypothesis (Tversky 1969, Rubinstein 1988).

Stage 1: The decision-maker looks for dominance. If the first alternative is at least as good as the second along both dimensions and better on at least one, then the first alternative is chosen.

Stage 2: If Stage 1 does not yield a verdict, the decision-maker compares each dimension separately, looking for similarities. If they perceive similarity in one dimension only, they prefer the alternative that is superior along the dissimilar dimension.

Stage 3: If neither Stage 1 nor Stage 2 yields a verdict, the choice is made using an unspecified different criterion.

There are several reasons that one might use this heuristic. First, because it checks for dominance, it avoids errors that might occur in an overall evaluation of each alternative in isolation. If such overall evaluation were imprecise, it would sometimes select an alternative that was slightly worse on both dimensions, in violation of dominance.

Second, the heuristic draws on readily available and easily evaluable information. Similarity appears to be among the features of objects and alternatives that are routinely

and automatically registered by the perceptual system (Engel and Wang 2011). People also appear to find it easier to evaluate differences than absolute magnitudes (Tversky and Kahneman 1979, 1983).

Third, the procedure capitalizes on the fact that intra-dimensional evaluation is relatively simple, because it involves comparisons between features of alternatives that are expressed in the same units (Tversky 1969). Subjects may also lack settled judgments about how to balance a loss on one dimension against a gain in another, which gives them reason to use a heuristic that side-steps such trade-offs (Tversky 1972, Brandstätter et al. 2006, Manzini and Mariotti 2007, Tserenjigmid 2015).

Notwithstanding these advantages, use of the similarity heuristic may yield choices that are a mere artefact of the choice procedure. It may also lead to violations of principles of rational choice. For example, it can lead to violations of transitivity of strict preference—the requirement that if a decision-maker has a strict preference for alternative *A* over *B*, and for *B* over *C*, then they must strictly prefer *A* to *C*. Suppose that *A* is worse than and similar to *B* along the first dimension, and better than and dissimilar to *B* along the second dimension. The similarity heuristic then leads to a preference for *A* over *B*. Further suppose that *B* is worse than and similar to *C* along the first dimension, and better than and dissimilar to *C* along the second dimension. The similarity heuristic then leads to a preference for *B* over *C*. Finally, suppose that *A* and *C* are dissimilar along both dimensions and that the first dimension is an important determinant of choice when alternatives differ substantially along it. Then, consistently with use of the similarity heuristic, the subject may prefer *C* to *A* in a pairwise comparison.

Following Tversky (1969), the role of similarity in choice has been studied in many experiments. The vast majority of these focus on choices between gambles (Lindman and

Lyons 1978, Budescu and Weiss 1987, Mellers et al. 1992, Leland 1994, Raynard 1995, Buschena and Zilberman 1995, 1999, Goldstone et al. 1997, Day and Loomes 2010, Loomes 2010, Regenwetter et al. 2011, Brandstätter and Gussmack 2013, Loomes and Pogrebna 2014), but some also examine the influence of similarity on other choices, including the trade-off between commuting time and wage (Mellers and Biagini 1994) and inter-temporal trade-offs (Rubinstein 2003). The predominant findings are: (i) in a substantial share (e.g., in Tversky's experiments, around one-third) of subjects, similar dimensions receive less weight in decision-making than dissimilar dimensions; and (ii) these subjects are prone to violating transitivity. (Because subjects are known not to choose deterministically, experiments standardly focus on Weak Stochastic Transitivity, which allows for a random error in the process of choice. In repeated pairwise choices between alternatives, this requires that if the probability of choosing *A* over *B* is greater than half and the probability of choosing *B* over *C* is greater than half, then the probability of choosing *A* over *C* must exceed one-half. In what follows, we also employ this conception of transitivity.)

This paper's intended contributions lie at the intersection of psychology and theories of distributive justice. First, we test for the use of the similarity heuristic in the little-studied moral domain. The vast majority of extant experiments on similarity-based decision-making involve self-regarding choices.¹ One may reasonably expect people to employ the same heuristics that they use in non-moral decision-making to make moral decisions. There are indeed a number of attempts to explain moral judgments as resulting from the application

¹ To our knowledge, the sole exceptions are one of Tversky's (1969) experiments, which involved pairwise choices between potential university applicants and Mellers' (1982) test of the influence of similarity on fairness judgments.

of decision principles whose use has been established primarily in non-moral domains, such as the principle that decisions are coded in relation to a baseline from which losses repel more than gains attract (see, e.g., Baron 1993, 1998; Kahneman 1994; Horowitz 1998; Sunstein 2005). However, such attempts have met with scepticism from philosophers, who have offered competing explanations of the judgments in question as the result of the application of attractive moral principles (see, e.g. Kamm 2007, pp. 422-49; Railton 2014). It is therefore worth establishing whether the similarity heuristic may determine judgment and choice in the domain of distributive justice and, in particular, whether it will yield choices that are contrary to all leading principles of justice. Interestingly, this project permits a new test of the use of the similarity heuristic. Most experiments infer the prevalence of similarity-based decision-making from violations of formal principles of rational choice (such as transitivity). While we do so too, we also use violations of substantive principles of distributive justice to diagnose its use.

Our second intended contribution is to the development of moral theory. Reaching reflective equilibrium in the area of distributive justice requires finding principles that explain and justify confidently held case judgments (Rawls 1999). We argue that judgments in cases that induce use of the similarity heuristic are not trustworthy and that they should therefore be excluded from the verdicts that principles of distributive justice should try to explain.

2. General idea of the experiment

How *should* one make trade-offs between the number of people one can save and the magnitude of the harm one can save them from? Standard theories of distributive justice that respect the Pareto principle range from utilitarianism (which requires maximizing the

sum-total of utility, or well-being, generated) to leximin (which requires maximizing the situation of the least-well-off). All such views except utilitarianism are willing to sacrifice some total utility for the sake of improving the lot of the worst off. And all such views except leximin are willing to accept some worsening in the situation of the worst off for the sake of a sufficiently large improvement in others' utility. This is true, for example, of forms of pluralist egalitarianism that care about both reducing inequality and improving total well-being (see Tungodden 2003). It is also true of the view known as prioritarianism, which does not care about inequality itself, but which gives some, non-infinite, extra weight to gains in utility that take place from a lower level (see Adler 2012).

How *do* subjects make these trade-offs? In line with the aversion to making trade-offs mentioned in the introduction, some simply avoid them. For example, in one study, Rodriguez-Miguez and Pinto-Prades (2002) report that 26% of subjects choose on the basis of one characteristic only: either they always save the greater number, even when saving the greater number does not maximize total utility, or they always save those facing the greatest harm, even when saving the better off would do far more good in aggregate. However, the predominant finding across many studies is that when subjects do make trade-offs, they tend to give substantial, though finite, extra weight to gains in well-being to the less well off (Nord and Johansen 2014).² Their moral preferences therefore generally align with the aforementioned pluralist egalitarian or prioritarian theories.

² Rodriguez-Miguez and Pinto-Prades (2002) report that this tendency disappears when subjects must choose between giving a large benefit to a single badly off individual and a small individual benefit to many badly off individuals. In such cases, contrary to all standard theories of distributive justice, subjects tend to offer the large, "concentrated" benefit to

To test the hypothesis that a substantial share of subjects would use the similarity heuristic, we proceeded as follows. We constructed pairs of alternatives that involved a trade-off between the magnitude of harm that people were saved from and the number of people saved so that: (i) these alternatives would appear similar along the “magnitude of harm averted” dimension but dissimilar along the “number of people saved” dimension; and (ii) choosing to save the more numerous group from the somewhat smaller harm would involve *helping the better off at a cost in total utility*. On these alternatives, we hypothesized, subjects who used the similarity heuristic would favour the more numerous better off, contrary to all aforementioned standard theories of distributive justice and contrary to the moral preferences generally evinced in surveys.

To establish subjects’ preferences when similarity could not determine choice, we also designed “wholly dissimilar” alternatives. Confronted with these alternatives, we conjectured, subjects would tend to help the less numerous worse off, both when this maximized total utility and when helping the worse off would come at some (modest) cost in total utility.

The conjectured switch between aiding the more numerous better off in choices between partly similar alternatives and aiding the less numerous worse off otherwise would

the single individual, even when “spreading” small benefits among many individuals would generate higher total utility and reduce inequality. Non-standard moral theories that permit such concentration of benefits are discussed in Temkin (2012) and Voorhoeve (2014). It is noteworthy that these non-standard theories cannot rationalize the choices to aid the better off at a cost in total utility in our experiment, in which each better off person stands to gain *less* than each worse off person.

be explicable neither in terms of standard theories of distributive justice nor by the use of the alternative heuristics that avoid trade-offs mentioned by Rodriguez-Miguez and Pinto-Prades (2002), viz. “always help the greater number” or “always help the worst off.” It would, however, be explained by use of the similarity heuristic.

3. Methods

We recruited 82 subjects (72% students, 28% non-students; 51% male, 49% female) from the subject pool of the Centre for Economic Learning and Social Evolution (ELSE) at University College London.

Subjects were sat in separate cubicles at individual computer screens. They were informed that they would be paid a flat fee of £13 (roughly USD 20 at the time) for participating in a 40-minute experiment on making choices in the use of health care resources and that a further £5 (USD 8) would be donated to a health care charity of their choice at the end of the experiment. (A full description of the introduction and questionnaire is available in [Appendix A1](#).)

Subjects were informed that they would face a series of choices between two interventions and asked which of the two the National Health Service should prioritise. Subjects were asked to suppose that the people affected were in their mid-thirties and in perfect health until recently, but that they now faced a health problem which diminished their well-being to the indicated level. If left untreated, these people would live the rest of a normal human lifespan with the indicated level of well-being; if treated, they would be returned to perfect health for the remainder of this lifespan. The measure of well-being used was the Health Utilities Index Mark III (Feeny et al. 1995, HUI Inc. 2008). This assigns 0 to death, 1 to perfect health, and a value in between to life in a state of impaired

functioning that is better than death.³ Subjects were told that it was developed by experts and were given a four-screen tutorial on its meaning. This included a picture of the scale, along with the representative valuation of eight conditions. It was explained that the values assigned to life in these conditions were determined by representative answers on surveys and that these values indicated the typical impact on well-being of a condition, with lower numbers representing lower well-being. Subjects were also informed that, on this scale, an increment of a given size always did a person just as much good, no matter from what level this increment took place.⁴

After this introduction, subjects were presented with four practice choices. The main experiment consisted of three “rounds” of going through sixteen choices in individually randomized order, for a total of forty-eight choices. (Every choice in the main experiment was therefore made three times, with, on average, fifteen other choices between repetitions.) After they had completed their choices, subjects were asked to offer a written explanation of five of their choices.

³ This index relies on the so-called “standard gamble” (Dolan 2001). If subjects respect the von Neumann-Morgenstern axioms, then it is a measure of von Neumann-Morgenstern utility.

⁴ Such stipulations notwithstanding, subjects may treat utility scales as if they have diminishing marginal prudential value (Greene and Baron 2001). If this were true of our subjects, this would make it even more difficult to achieve the hypothesized preference for aiding the better off at a cost in total utility, and so a preference of this kind would be even more strongly indicative of the use of similarity-based decision-making.

Alternatives were displayed as in Figure 1. (The placement of alternatives on the right-hand and left-hand side of the screen was randomized.) In the figure, the solid vertical line and number to the left alongside the top of this line represent the group's health status if untreated (0.95 for alternative *A* and 0.91 for *B*). The dotted line represents the health that would be restored by treatment. A box attached to the top of the solid line contained the number of people in that condition that one can treat (48 for *A* and 27 for *B*). Subjects were told that they could treat only one of the two groups. They did so by clicking on the box with the number of patients in that group and moving it all the way up to full health.

4. Results and discussion

As a test of basic comprehension and attentiveness, we included choice in which one alternative dominated another (*N* versus *O* in Table 1). Seventy-six subjects (92.7%) passed this test with flying colours, selecting the dominant alternative three times out of three. A further three subjects (3.7%) chose the dominant alternative two times out of three. Another three subjects (3.7%) chose the dominated alternative two or more times out of three. We exclude the latter three subjects from the following analysis. (This exclusion makes no substantial difference to our conclusions.)

4.1 Alternatives that are similar in terms of health gain

Consider the choice between *A* and *B* in Figure 1. We conjectured that in this choice, the alternatives would appear similar along the "health-related utility gained" dimension, but dissimilar along the "number of people helped" dimension. Use of the similarity heuristic would therefore yield a preference for *A*, despite the fact that *B* would both aid the worst off and yield greater total utility. As listed in Table 1, we constructed a further three

alternatives, *C*, *D*, and *E*, each of which was designed to appear similar to its immediate predecessor along the health-related utility gain dimension and each of which yields greater total utility than its predecessor.

As Table 2 reveals, helping the better off at a cost in total utility is indeed common in the pairwise choices between these alternatives. For *A* versus *B*, *B* versus *C*, and *C* versus *D*, more than half of all subjects favour the better off at least 2 out of the 3 times they were presented with the choice. This implies that, if decisions in these pairwise comparisons were taken by majority voting, our group of subjects would choose contrary to every leading theory of distributive justice.

The exception is *D* versus *E*, in which, at 38%, the preference for aiding the better off is less common than in the other choices between adjacent alternatives in the *A* to *E* sequence. (Statistical tests reported in [Appendix A2](#), Table A2.2 confirm that *D* versus *E* stands apart.) Our explanation for this is that, as we move through this sequence, the absolute difference in (and the ratio of) the number of people saved shrinks to the point that some subjects would perceive *D* and *E* as similar along *both* dimensions. For such subjects, the similarity heuristic does not mandate a choice at Stage 2. Instead, it moves to Stage 3, at which we conjectured that subjects would display preferences in line with standard theories of distributive justice, all of which mandate aiding the worse off in this choice.

Because the gap in health gain between non-adjacent alternatives in the set *A* to *E* is larger than between adjacent alternatives, the former are more likely to look dissimilar along both dimensions. A key prediction of our hypothesis is therefore that subjects would be less likely to aid the better off (at a cost in total utility) in pairwise choices between non-adjacent alternatives than in choices between adjacent alternatives. As Figure 2 shows, this

is indeed what occurred. To establish whether this difference in the rate of aiding the better off is statistically significant, we consider the comparisons listed in Table 3. The underlined numbers in the top-right corner of each comparison indicate the share of subjects who engage in the predicted switching from aiding the better off in a choice between adjacent alternatives to aiding the worse off in a choice between non-adjacent alternatives. (For example, 22.8% of subjects both aid the better off in *A* versus *B* and aid the worse off in *A* versus *C*.) There is no comparable shift in the opposite direction. (For example, only 6.3% of subjects aid the worse off in *A* versus *B* and the better off in *A* versus *C*.) We use McNemar's exact test to calculate the probability of these results given the null hypothesis that the answers to two different pairwise choices are random draws from the same binomial distribution. (For discussion of this test, see [Appendix A2](#).) The grey box in the middle of each comparison in Table 3 reports the results. The differences between adjacent and non-adjacent choices are significant at at least the 5% level throughout.

This switch from favouring the better off in a choice between adjacent alternatives to favouring the worse off in a choice between non-adjacent alternatives should make it more likely that subjects will display intransitive preferences. For example, subjects using the similarity heuristic would favour *A* over *B* and *B* over *C*. But if they also found *A* wholly dissimilar to *C*, they would, in line with standard theories of distributive justice, prefer *C* to *A*, violating transitivity.

To assess this prediction, we divide our subjects into three groups: those who do not make intransitive choices (59.5%), those who make intransitive choices in a manner explicable in terms of use of the similarity heuristic (35.4%), and those who make intransitive choices that are not so explicable (5.1%). As Table 4 shows, we can confidently reject the hypothesis that the latter two groups are equally large.

Intransitivities of the kind induced by similarity-based decision-making are also manifest at the group level. As Table 2 reveals, pairwise majority voting yields a group preference for A over B , B over C , and C over D , but it also a preference for C over A , D over A , and B over D , yielding three intransitive cycles.

4.2 Wholly dissimilar alternatives

Despite the larger gap in terms of health gain, some subjects might still have found some non-adjacent alternatives in the set A through E similar along the health-related utility gain dimension. After all, the difference in health gain between, say, A and C is only 0.08. We therefore constructed further choices between wholly dissimilar alternatives, all listed in Table 1. Each of the pairwise choices R versus S , T versus U , and V versus W involves a stark choice between helping a substantially smaller, substantially worse off group and helping a much larger, much better off group. In all cases, helping the worse off yields somewhat greater utility. Figure 3 reveals that, as predicted, aiding the better off at a cost in total utility is indeed much more frequent in choices among partly similar alternatives than in choices among wholly dissimilar alternatives. (Our analysis in Table A2.3 in [Appendix 2](#), confirms that this difference is statistically significant at at least the 5% level.)

To test our conjecture that a substantial number of subjects will both aid the better off at cost in total utility in choices between partly similar alternatives and aid the worst off at a cost in total utility in choices between wholly dissimilar alternatives, we constructed G versus F and Q versus P . Table 5 shows that our evidence supports this conjecture. For example, a striking 49.4% of all subjects shift from aiding the better off in the choice between A and B to aiding the worse off in the choice between G and F . (Only 5.1% switch in

the opposite direction.) These differences are statistically significant at the 1% level for all pairs, save one, where this switch is significant only at the 10% level.

4.3 Subjects' decision rules

We shall now examine individual-level data. We start by matching individuals with the decision rule that best represents their choices. In doing so, we must note that the similarity heuristic is consistent with a wider range of choices than the other decision rules under examination, because it allows for individual-level variability in perceptions of similarity. This flexibility may give the similarity heuristic an “unfair advantage” over other decision rules. We attenuate this problem as follows. We first report an analysis which allows some, albeit limited, variability in individuals' perceptions of similarity. We then consider how robust our findings are by imposing the same perceptions of similarity on all subjects.

Our first test permits only the following two types of perceptions of similarity:

Adjacent Only: All adjacent alternatives in the *A* through *E* sequence, and only these alternatives, are similar along the gain dimension;

Two Steps Only: All alternatives that are no more than one step apart in the *A* through *E* sequence, and only these alternatives, are similar along the gain dimension.

Moreover, we do not consider data from *D* versus *E*, on which the similarity heuristic would have an unfair advantage because it is consistent with either choice, since *D* and *E* may be regarded as similar along the gain dimension only (in which case it predicts that *D* is chosen), or along both dimensions (in which case it predicts that *E* is chosen). We also do not consider *N* versus *O*, since this was a mere basic comprehension test. Since each subject confronted each of the remaining fourteen comparisons three times, this yields 42 data

points for each individual. We then assign each individual to the decision rule that gets the largest share of these choices right. The second column of Table 6 displays the results. It indicates that the similarity heuristic is the most common decision rule, with some form of it being the best description of 41.8% of the population. It also reveals that almost all of those who use the similarity heuristic are prepared to sacrifice total utility for the sake of the worse off when alternatives are wholly dissimilar. The second-most popular decision rule (the uniquely best match for 35.4% and tied for best for a further 5.1%) is to always help the worst off. A small minority (12.7%) is best described as always saving the greater number; an even smaller minority (5.1%, or 10.1% if one counts ties) is best described as maximizing total utility.

The third column of Table 6 indicates how well these decision rules fit the choices of the subjects matched with them. The similarity heuristic fits its matched population reasonably well, with a “success rate” of 78.5%. The final column reveals that this heuristic adds substantially to our ability to predict the choices of the subjects matched with it—on average, if we could not use this heuristic to describe their choices, our success rate at describing them would drop by 11.6%.

As a robustness check, we also considered a version of the similarity heuristic that imposes the uniform perception of (dis)similarity inherent in the aforementioned Adjacent Only rule on all subjects. This is a demanding test of this heuristic, since some diversity in individual perceptions of similarity is to be expected, and does not imply that individuals do not use the heuristic. As detailed in [Appendix 2](#), Table A2.8, this imposition lowers the share of subjects whose behaviour best matches the similarity heuristic to 36.7%, placing it on a par with “always aid the worse off”. Nonetheless, the results of this test support the idea that the similarity heuristic is used by around two-fifths of subjects.

Further evidence can be gleaned from subjects' written explanations of five of their choices, which they completed at the end of the experiment. We pigeonholed each of their answers using one of the five categories listed in the first column of Table 7. To illustrate this categorization, consider the following examples of subjects' explanations of their choices in *A* versus *B*. (Subjects' complete answers and our categorization can be accessed in [Appendix A3](#).)

S44 offered the following explanation of their preference for *A*:

"Because it helped 21 more people & there was only 0.04 difference in severity of problem."

We categorized this answer as displaying evidence of use of the similarity heuristic.

S47 explained their preference for *B* as follows:

"48 people are almost fine. 27 are worse off; so they should be helped."

We categorized this answer as expressing a special concern for the worse off.

S60 offered the following explanation for their preference for *A*:

"The more number of people to treat [sic]."

We categorized this answer as indicating adherence to a rule requiring saving the greater number.

S41 explained their preference for *B* over *A* as follows:

"the total gain is more because $27 * 0.09 > 0.05 * 48$."

This answer was categorized as expressing adherence to the rule of maximizing total utility.

S81, whose choices expressed a preference for *A*, wrote:

"other people in reasonable health."

This answer was categorized as not rationalizing the choice in question, since it is too terse to serve as a justification. (Other reasons for placing responses in this category were offering reasons that justified choices that differed from the subjects' actual choices, or not answering the question.)

Several results are worth highlighting. First, for choices between alternatives that are similar along one dimension, the most frequently offered explanation involves this similarity. Second, in wholly dissimilar choices, concern for the worse off predominates as an explanation. Third, as the final two columns of Table 7 make clear, this increase in attention to the worse off is almost entirely due to subjects who switch from appealing to similarity to justify their choice to appealing to the fate of the worse off. Finally, other rationales, including the aim of maximizing total utility, are infrequently invoked.

Our individual-level data therefore helps assess a hypothesis raised by a number of commentators, which is that a substantial share of subjects aim to maximize total utility throughout, but simply make errors in estimating the alternative with the highest total utility in pairwise choices between adjacent alternatives in the *A* through *E* sequence. These errors, so this hypothesis goes, are committed because the total is difficult to calculate and the difference in total utility between the alternatives is small.⁵

We note that this “people are error-prone utilitarians” hypothesis is compatible with the idea that subjects use the similarity heuristic when they have difficulty engaging in the

⁵ This hypothesis was raised by Antonio Cabrales and Joseph Millum as well as by seminar audiences. It gains indirect support from the finding in Arieli et al. (2011) that subjects were more likely to engage in separate evaluation of the probability and prize dimensions of gambles when the expected monetary value of the gamble was difficult to compute.

“holistic” evaluation of alternatives. For it is consistent with the idea that subjects use this heuristic to estimate which alternative maximizes total utility when calculating this total is demanding. Nonetheless, if correct, it would conflict with our idea, mentioned in the introduction, that in cases of the kind under consideration, many subjects do not yet have a fully articulated theory of distributive justice which they are trying to apply. Our findings, however, offer very little support for the “people are error-prone utilitarians” hypothesis. As we have seen, utilitarianism best fits only a handful of people’s choices (see Table 6). This is because in *G* versus *F* and *Q* versus *P*, the vast majority chose to aid the worse off at a cost in total utility. (It is noteworthy that these are choices in which total utility was relatively easy to calculate.) Finally, as the final column in Table 7 reveals, in only 2% of all cases in which a subject invoked similarity as a rationale for aiding the better off did they also invoke a utilitarian rationale for their choices between wholly dissimilar alternatives.

In sum, subjects’ accounts of their choices confirm the conclusions we drew from our analysis of individual-level choice data. Both types of evidence indicate that, with roughly two-fifths employing it, the similarity heuristic is the most commonly used decision rule (closely followed by special concern for the worse off). Moreover, both subjects’ choices and their proffered rationales indicate that the vast majority of individuals who employ the similarity heuristic choose on the basis of concern for the worse off when faced with wholly dissimilar alternatives.

5. Empirical and normative conclusions

We have examined how people choose when they must either save a larger number of people from a smaller harm, or, instead, save a smaller number of people from a greater harm. We have documented a remarkable shift in subjects’ decisions. In choices between

alternatives that appear similar only along the “magnitude of harm prevented” dimension, a majority of subjects help the more numerous better off at a cost in total utility. By contrast, in choices between wholly dissimilar alternatives, a vast majority of subjects help the less numerous worse off, even when this comes at a cost in total utility. This shift leads to violations of transitivity at the individual and collective level. We have argued that these patterns of choice are best explained by widespread use of the similarity heuristic; indeed, both individual-level choice data and subjects’ written accounts of their choices indicate that somewhat in excess of 40% of subjects employ this heuristic.

In our experiment, these subjects’ choices do not express a consistent set of values. We conclude from this that we should not take all of subjects’ choices in our experiment to accurately represent their considered moral preferences. Confronted with our experimental findings, a reflective subject displaying the similarity-induced pattern of choice should, we believe, revise some of their preferences. In carrying out such revision, we submit that such a subject should consider that the similarity-induced choices are most likely the result of not fully engaging in the difficult trade-off in question. They should also consider that these choices are, taken separately, at odds with every standard theory of distributive justice and, taken as a set, are often at odds with formal requirements of rational choice. To us, these facts indicate that their similarity-induced choices are mere artefacts of the choice situation and are therefore untrustworthy. Indeed, it seems to us likely that when deciding between alternatives that are similar only in terms of the harm from which individuals can be saved, similarity-based reasoning leads subjects to systematically underweight this harm’s importance.

We conclude that decisions between pairs of alternatives that are similar along precisely one dimension are unsuitable for the projects mentioned in the introduction,

namely surveying the public to uncover its moral preferences and using data from thought experiments to test and develop theories of distributive justice. Our recommendation to both social scientists creating survey instruments and philosophers designing thought experiments is therefore to construct choices that do not induce subjects to use the similarity heuristic.⁶

Acknowledgements

This paper was presented at Copenhagen University, the Experimental Philosophy Conference at Newcastle University, Fudan University, Harvard University, the Institute of Philosophy in London, LSE, Oxford University, Universidade Nova, the University of Maryland at College Park, Uppsala University, the U.S. National Institutes of Health, and Warwick University. We thank our audiences, Ken Binmore, Luc Bovens, Antonio Cabrales, Ipek Gençsü and Joseph Millum for comments and UCL's Centre for Economic Learning and Social Evolution (ELSE) for the use of its laboratory. This research was supported by the British Academy through grant SG 45949 and by the British Arts and Humanities Research Council through grant AH/J006033/1. The opinions expressed are the view of the authors only. They do not represent any position or policy of the U.S. National Institutes of Health, the Public Health Service, or the Department of Health and Human Services.

⁶ Philosophers, at least, have not always followed this dictum. For examples of thought experiments that arguably induce use of the similarity heuristic, see Quinn (1990), Rachels (1998) and Temkin (2012). For a critique of these thought experiments, see Voorhoeve and Binmore (2006) and Voorhoeve (2008, 2013).

References

- Adler, M. 2012. *Well-Being and Fair Distribution: Beyond Cost-Benefit Analysis*. Oxford: Oxford University Press.
- Arieli, A., Y. Ben-Ami, and A. Rubinstein. 2011. Tracking decision-makers under uncertainty. *American Economic Journal: Microeconomics*, 3: 68—76. DOI: 10.1257/mic.3.4.68
- Baron, J. 1993. *Morality and Rational Choice*. Dordrecht: Springer.
- Baron, J. 1998. *Judgment misguided: Intuition and error in public decision making*. Oxford: Oxford University Press.
- Brandstätter, E. and E. Gussmack. 2013. The cognitive processes underlying risky choice. *Journal of Behavioral Decision Making*, 26: 185—197. DOI: 10.1002/bdm.1752
- Brandstätter, E., G. Gigerenzer, and R. Hertwig. 2006. The priority heuristic: Making choices without trade-offs. *Psychological Review* 113: 409—32. DOI: 10.1037/0033-295X.113.2.409
- Budescu, D. and W. Weiss. 1987. Reflection of transitive and intransitive preferences: A test of prospect theory. *Organizational Behaviour and Human Decision Processes* 39: 184—202. DOI: 10.1016/0749-5978(87)90037-9
- Buschena, D. and D. Zilberman. 1995. Performance of the similarity hypothesis relative to existing models of risky choice. *Journal of Risk and Uncertainty* 11: 233—62. DOI: 10.1007/BF01207788
- Buschena, D. and D. Zilberman. 1999. Testing the effects of similarity on risky choice: Implications for violations of expected utility. *Theory and Decision* 46: 251—76. DOI: 10.1023/A:1005066504527

- Daniels, Norman. 2013. Reflective equilibrium. *The Stanford Encyclopedia of Philosophy* (Winter 2013 Edition), Edward N. Zalta (ed.), URL = <http://plato.stanford.edu/archives/win2013/entries/reflective-equilibrium/>.
- Day, B. and G. Loomes. 2010. Conflicting violations of transitivity and where they may lead us. *Theory and Decision* 68: 233–242. DOI: 10.1007/s11238-009-9139-1
- Dhar, R., S. M. Nowlis, and S. J. Sherma. 1999. Comparison effects on preference construction. *Journal of Consumer Research* 26: 293—306. DOI: 10.1086/209564
- Dolan, P. 2001. Output measures and valuation in health. In *Economic Evaluation in Health Care*, edited by Michael Drummond and Alistair McGuire, 46—67. Oxford: Oxford University Press.
- Drechsler, M., K. Katsikopoulos, G. Gigerenzer. 2014. Axiomatizing bounded rationality: the priority heuristic. *Theory and Decision* 77: 183—96. DOI: 10.1007/s11238-013-9393-0
- Engel, T. and X.-J. Wang. 2011. Same or different? A neural circuit mechanism of similarity-based pattern match decision making. *The Journal of Neuroscience* 31(19): 6982–96. DOI:10.1523/JNEUROSCI.6150-10.2011
- Feeny, D., W. Furlong, M. Boyle, and G. W. Torrance. 1995. Multi-attribute health status classification systems: Health Utilities Index. *Pharmacoeconomics* 7: 490—502. DOI: 10.2165/00019053-199507060-00004
- Gaertner, W. and E. Schokkaert. 2012. *Empirical Social Choice: Questionnaire-Experimental Studies on Distributive Justice*. Cambridge: Cambridge University Press.
- Gigerenzer, G., P. Todd, and the ABC Research Group. 2000. *Simple Heuristics that Make Us Smart*. Oxford: Oxford University Press.

- Goldstone, R., D. Medin, and J. Halberstadt. 1997. Similarity in context. *Journal of Memory and Cognition* 25: 237—55. DOI: 10.3758/BF03201115
- Greene, J. and J. Baron. 2001. Intuitions about declining marginal utility. *Journal of Behavioral Decision Making* 14: 243—55. DOI: 10.1002/bdm.375
- Horowitz, T. 1998. Philosophical intuitions and psychological theory. *Ethics* 108: 367—85. DOI: 10.1086/233809
- HUI Inc. 2008. The Health Utilities Index Mark 3. <http://www.healthutilities.com/> [Last accessed October 13, 2016].
- Kahneman, D. 1994. The Cognitive Psychology of Consequences and Moral Intuition. Delivered as a Tanner Lecture on Moral Values, unpublished manuscript.
- Kamm, F.M. 2007. *Intricate Ethics*. Oxford: Oxford University Press.
- Köszegi, B. and A. Szeidl. 2013. A model of focusing in economic choice. *Quarterly Journal of Economics* 128(1): 53—104. DOI: 10.1093/qje/qjs049
- Leland, J. 1994. Generalized similarity judgments: An alternative explanation for choice anomalies. *Journal of Risk and Uncertainty* 9: 151—72. DOI: 10.1007/BF01064183
- Lindman, H. and J. Lyons. 1978. Stimulus complexity and choice inconsistency among gambles. *Organizational Behavior and Human Performance* 21: 146—59. DOI: 10.1016/0030-5073(78)90046-6
- Loomes, G. 2010. Modeling choice and valuation in decision experiments. *Psychological Review* 117(3): 902—24. DOI: 10.1037/a0019807
- Loomes, G. and G. Pogrebna. 2014. Testing for independence while allowing for probabilistic choice. *Journal of Risk and Uncertainty* 49: 189—211. DOI: 10.1007/s11166-014-9205-0

- Manzini, P. and M. Mariotti. 2007. Sequentially rationalizable choice. *American Economic Review* 97 (5): 1824—39. DOI: 10.1257/aer.97.5.1824
- Mellers, B. 1982. Equity judgment: A revision of Aristotelian views. *Journal of Experimental Psychology: General* 111: 242—70. DOI: 10.1037/0096-3445.111.2.242
- Mellers, B., S. Chang, M. Birnbaum and L. Ordonez. 1992. Preferences, prices, and ratings in risky decision making. *Journal of Experimental Psychology: Human Perception and Performance* 18: 347—61. DOI: 10.1037/0096-1523.18.2.347
- Mellers, B. and K. Biagini. 1994. Similarity and choice. *Psychological Review* 101: 505—18. DOI: 10.1037/0033-295X.101.3.505
- Nord, E. and R. Johannsen. 2014. Concerns for severity in priority setting in health care: A review of trade-off data in preference studies and implications for societal willingness to pay for a QALY. *Health Policy* 116: 281—8. DOI: 10.1016/j.healthpol.2014.02.009
- Quinn, W. 1990. The puzzle of the self-torturer. *Philosophical Studies* 59: 79—90. DOI: 10.1007/BF00368392
- Rachels, S. 1998. Counterexamples to the transitivity of better than. *Australasian Journal of Philosophy* 76: 71—83. DOI: 10.1080/00048409812348201
- Railton, P. 2014. The Affective Dog and Its Rational Tale: Intuition and Attunement. *Ethics* 124: 813—59. DOI: 10.1086/675876
- Raynard, B. 1995. Reversals of preference between compound and simple risks: The role of editing heuristics. *Journal of Risk and Uncertainty* 11: 159—75.
- Rawls, J. 1999. *A Theory of Justice, revised, 2nd edition*. Oxford: Oxford University Press.
- Regenwetter, M., J. Dana, C. Davis-Stober. 2011. Transitivity of preferences. *Psychological Review* 118: 42—56. DOI: 10.1037/a0021150

- Rodriguez-Miguez, E. and J.-L. Pinto-Prades. 2002. Measuring the social importance of concentration or dispersion of individual health benefits. *Health Economics* 11: 43—53. DOI: 10.1002/hec.643
- Rubinstein, A. 1988. Similarity and decision-making under risk (Is there a utility theory resolution to the Allais Paradox?) *Journal of Economic Theory* 46: 145—53. DOI: 10.1016/0022-0531(88)90154-8
- Rubinstein, A. 2003. Economics and psychology? The case of hyperbolic discounting. *International Economic Review* 44: 1207—16. DOI: 10.1111/1468-2354.t01-1-00106
- Sunstein, C. 2005. Moral heuristics. *Behavioural and Brain Sciences* 28: 531—73. DOI: 10.1017/S0140525X05000099
- Temkin, L. 2012. *Rethinking the Good: Moral Ideals and the Nature of Practical Reasoning*. Oxford: Oxford University Press.
- Tserenjigmid, G. 2015. Theory of decisions by intra-dimensional comparisons. *Journal of Economic Theory* 159: 326—38. DOI: 10.1016/j.jet.2015.07.001
- Tungodden, B. 2003. The value of equality. *Economics and Philosophy* 19: 1—44. DOI: 10.1017/S0266267103001007
- Tversky, A. 1969. Intransitivity of preferences. *Psychological Review* 76: 31—48. DOI: 10.1037/h0026750
- Tversky, A. 1972. Elimination by aspects: A theory of choice. *Psychological Review* 79: 281—99. DOI: 10.1037/h0032955
- Tversky, A. and D. Kahneman. 1974. Judgment under uncertainty: Heuristics and biases. *Science* 185: 1124—31. DOI: 10.1126/science.185.4157.1124
- Tversky, A. and D. Kahneman. 1979. Prospect theory. *Econometrica* 47: 263—91. DOI: 10.2307/1914185

- Tversky, A. and D. Kahneman. 1983. Extensional vs. intuitive reasoning: The conjunction fallacy in probability judgment. *Psychological Review* 90: 293—315. DOI: 10.1037/0033-295X.90.4.293
- Ubel, P.A., G. Loewenstein, D. Scanlon, and M. Kamlet. 1996. Individual utilities are inconsistent with rationing choices: A partial explanation of why Oregon's cost-effectiveness list failed. *Medical Decision-Making* 16: 108—16. DOI: 10.1177/0272989X9601600202
- Voorhoeve, A. and K. Binmore. 2006. Transitivity, the sorites paradox, and similarity-based decision-making. *Erkenntnis* 64: 101—14. DOI: 10.1007/s10670-005-2373-1
- Voorhoeve, A. 2008. Heuristics and biases in a purported counterexample to the acyclicity of 'better than'. *Politics, Philosophy and Economics* 7: 285—99. DOI: 10.1177/1470594X08092104
- Voorhoeve, A. 2013. Vaulting intuition: Temkin's critique of transitivity. *Economics and Philosophy* 29: 409—25. DOI: 10.1017/S0266267113000321
- Voorhoeve, A. 2014. How should we aggregate competing claims? *Ethics* 125: 64-87. DOI: 10.1086/677022

Figure 1. A choice between alternatives *A* and *B*.

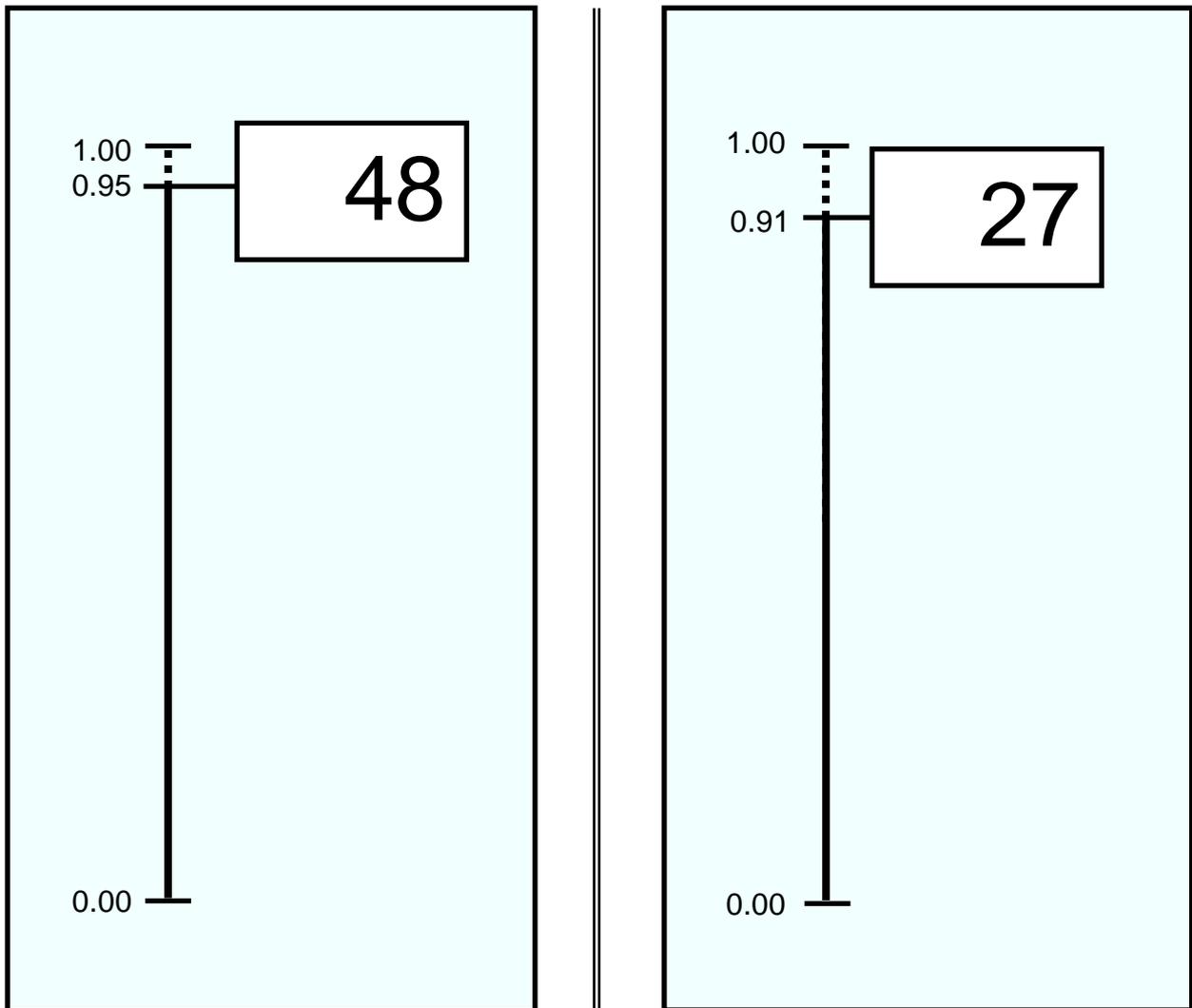


Table 1. All alternatives

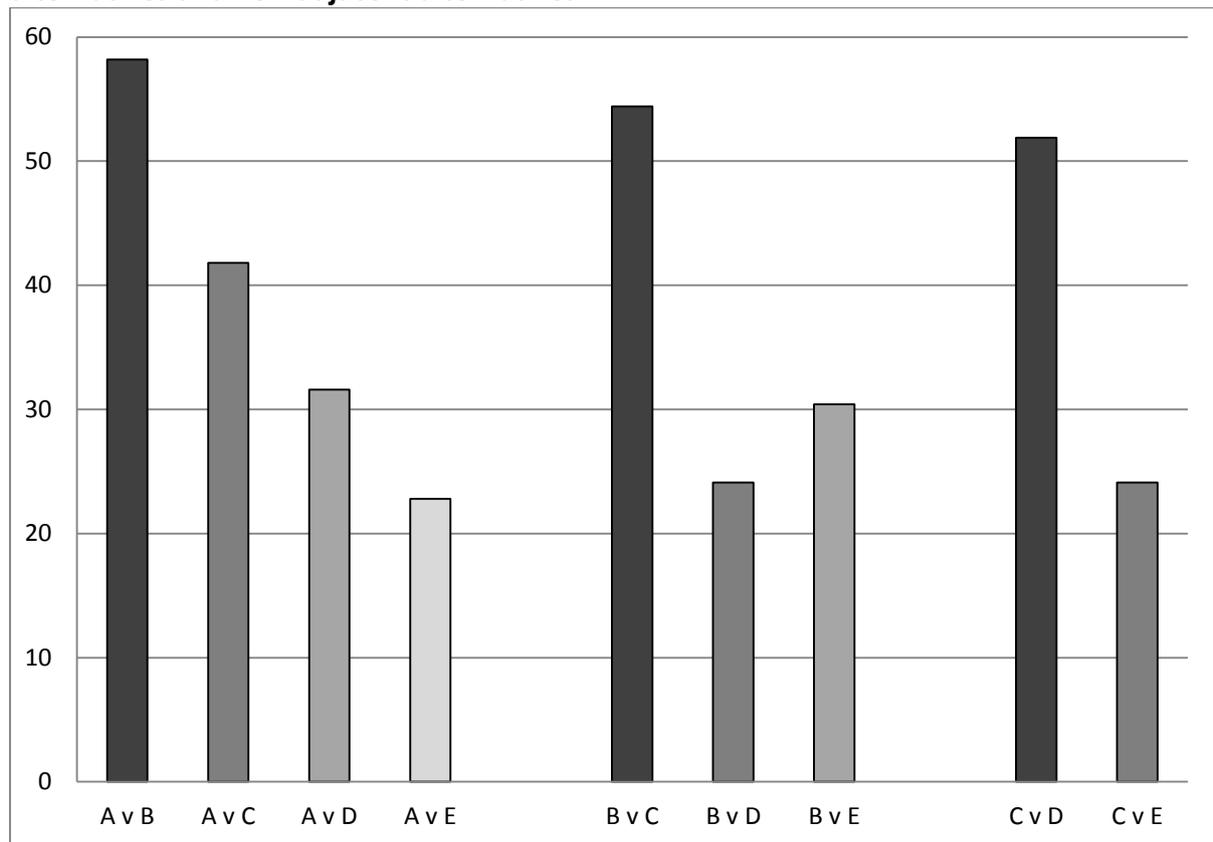
Alternatives		Number of patients	Individual utility without treatment	Individual gain through treatment	Sum of utilities generated
Similar to adjacent alternative along utility gain dimension; aiding better off lowers total utility	<i>A</i>	48	0.95	0.05	2.40
	<i>B</i>	27	0.91	0.09	2.43
	<i>C</i>	19	0.87	0.13	2.47
	<i>D</i>	15	0.83	0.17	2.55
	<i>E</i>	12	0.78	0.22	2.64
Note: Subjects choose between all possible pairings from the set <i>A</i> through <i>E</i> .					
Wholly dissimilar; aiding better off lowers total utility	<i>R</i>	10	0.31	0.69	6.90
	<i>S</i>	20	0.70	0.30	6.00
	<i>T</i>	5	0.26	0.74	3.70
	<i>U</i>	20	0.84	0.16	3.20
	<i>V</i>	5	0.50	0.50	2.50
	<i>W</i>	21	0.90	0.10	2.10
	Wholly dissimilar; aiding better off raises total utility	<i>F</i>	10	0.05	0.95
<i>G</i>		50	0.80	0.20	10.0
<i>P</i>		1	0.10	0.90	0.90
<i>Q</i>		4	0.75	0.25	1.00
Basic comprehension	<i>N</i>	31	0.14	0.86	26.70
	<i>O</i>	14	0.35	0.65	9.10

Table 2. Subjects expressing a preference for aiding the better off, in percent

Choices		Percentage of individuals ($n = 79$) expressing preference for better off (≥ 2 times out of 3)
Similar in terms of health gain; aiding better off lowers total utility	$A \vee B$	58.2
	$B \vee C$	54.4
	$C \vee D$	51.9
	$D \vee E$	38.0
Possible dissimilarity along both dimensions; aiding better off lowers total utility	$A \vee C$	41.8
	$A \vee D$	31.6
	$B \vee D$	24.1
	$B \vee E$	30.4
	$C \vee E$	24.1
Wholly dissimilar; aiding better off lowers total utility	$A \vee E$	22.8
	$S \vee R$	12.7
	$U \vee T$	13.9
	$W \vee V$	13.9
Wholly dissimilar; aiding better off raises total utility	$G \vee F$	13.9
	$Q \vee P$	24.1

Note: $n = 79$. In the pairwise choices in the second column, the alternative which involves aiding the better off is always listed first. Note also that one can also consider the percentage of all *choices* (rather than subjects) that favour the better off. The resulting pattern is very similar; see [Appendix A2](#), Table A2.1.

Figure 2. Number of subjects (in percent) aiding the better off in choices between adjacent alternatives and non-adjacent alternatives.



Note: $n = 79$. Dark bars indicate choices between adjacent alternatives, which are more likely to be perceived as similar in terms of health gain only. Lighter bars indicate choices between non-adjacent alternatives; these are more likely to be perceived as wholly dissimilar. They are lighter the further apart the alternatives are. Aiding the better off (at a cost in total utility) is much more frequent among adjacent alternatives than among non-adjacent alternatives.

Table 3. Comparison of choices between adjacent alternatives with choices between nonadjacent alternatives.

Choices between nonadjacent (more likely to be perceived as wholly dissimilar) alternatives

		<i>A v C</i>		<i>A v D</i>		<i>A v E</i>	
		Better off	Worse off	Better off	Worse off	Better off	Worse off
<i>A v B</i>	Better off	35.4	<u>22.8</u>	27.8	<u>30.4</u>	21.5	<u>36.7</u>
	Worse off	6.3	35.4	3.8	38.0	1.3	40.5
		<i>B v D</i>		<i>B v E</i>			
		Better off	Worse off	Better off	Worse off		
Choices between adjacent (partly similar) alternatives	Better off	21.5	<u>32.9</u>	26.6	<u>27.8</u>		
	Worse off	2.5	43.0	3.8	41.8		
		<i>C v E</i>					
		Better off	Worse off				
<i>C v D</i>	Better off	24.1	<u>27.8</u>				
	Worse off	0.0	48.1				

Note: $n = 79$. Numbers in the comparisons of distributions across (aiding the better off at a cost in total utility, aiding the worse off at a gain in total utility) are percentages of the total population. Underlined numbers represent the predicted shift from aiding the better off when choosing between adjacent alternatives to aiding the worse off when choosing among non-adjacent alternatives. Numbers in the grey boxes give the probability of obtaining the observed results if the answers come from the same distribution, using McNemar's exact test. We can reject the hypothesis that choices between adjacent alternatives and choices between nonadjacent alternatives come from the same distribution.

** = 5% confidence level.

*** = 1% confidence level.

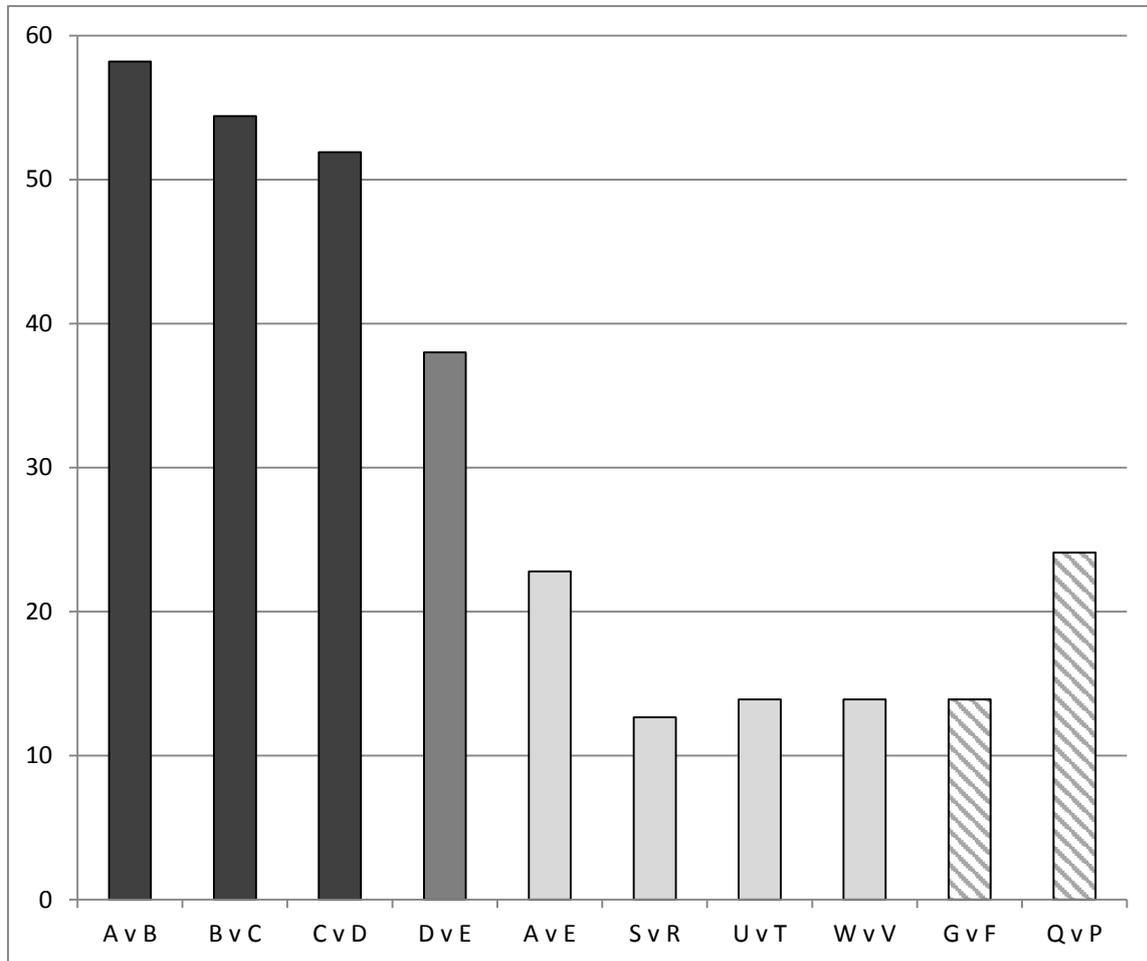
Table 4. Types of intransitivities observed.

	Proportion of population, in %	Probability no difference in direction of intransitivity
No intransitivities	59.5	
Intransitivity explicable by similarity	35.4	0.000***
Intransitivity not explicable by similarity	5.1	

Note: $n = 79$. The final column lists the probability p of obtaining the observed proportions if the probability q of each type of intransitivity were the same. Note that this can be done for any value of q between 0 and 0.5. Therefore, the q that maximises p is chosen. The probability is calculated using equation A2.2 in [Appendix 2](#). In line with our prediction, we can confidently reject the hypothesis that both types of intransitivity are equally likely.

*** = 1% confidence level.

Figure 3. Comparison of the rate of preference for the better off (in percent) in choices between alternatives that are similar in terms of health gain with this rate for choices among wholly dissimilar alternatives.



Note: $n = 79$. Darker bars indicate choices between alternatives that are more likely to involve alternatives that are similar in terms of health gain only; light bars indicate choices between wholly dissimilar alternatives. Bars with an even colouring indicate choices in which aiding the better off decreases total utility. Patterned bars indicate choices in which aiding the better off increases total utility.

Two conclusions are apparent. First, giving priority to the better off at a cost in total utility is much more frequent in choices among alternatives that are similar in terms of health gain. ([Appendix 2](#), Table A2.3 confirms that this difference is statistically significant.) Second, in choices between wholly dissimilar alternatives ($G v F$ and $Q v P$), the vast majority of subjects is prepared to sacrifice total utility for the sake of the worst off.

Table 5. The shift from aiding the better off to aiding the worse off.

		Wholly dissimilar choices; aiding better off raises total utility			
		<i>G v F</i>		<i>Q v P</i>	
		Better off	Worse off	Better off	Worse off
<i>A v B</i>	Better off	8.9	<u>49.4</u>	13.9	<u>44.3</u>
		0.000***		0.000***	
	Worse off	5.1	36.7	10.1	31.6
<i>B v C</i>	Better off	10.1	<u>44.3</u>	13.9	<u>40.5</u>
		0.000***		0.000***	
	Worse off	3.8	41.8	10.1	35.4
<i>C v D</i>	Better off	12.7	<u>39.2</u>	13.9	<u>38.0</u>
		0.000***		0.000***	
	Worse off	1.3	46.8	10.1	38.0
<i>D v E</i>	Better off	7.6	<u>30.4</u>	8.9	<u>29.1</u>
		0.001***		0.090*	
	Worse off	6.3	55.7	15.2	46.8
Total		13.9	87.1	24.1	75.9

Similarity in terms of health gain; aiding better off lowers total utility

Note: $n = 79$. Numbers in the comparisons of distributions across (aiding the better off at a cost in total utility, aiding the worse off at a cost in total utility) give percentages of the total population. Underlined numbers represent the predicted shift. Numbers in the grey boxes give the probability of obtaining the observed results if the answers come from the same distribution according to McNemar's exact test. A large share of subjects switch from aiding the better off at a cost in total utility when choosing between similar alternatives to aiding the worse off at a cost in total utility when choosing between wholly dissimilar alternatives.

* = 10% confidence level.

*** = 1% confidence level.

Table 6. Matching subjects with decision rules

Rule	Share (%)	Fit (%)	Fit premium (%)
Similarity heuristic	41.8	78.5	11.6
When no similarity:			
Worst off	40.5	78.9	
Total utility	1.3	66.7	
Worse off	35.4	88.8	8.3
Greater number	12.7	80.0	23.1
Total utility	5.1	79.2	4.8
Worse off/total utility (tie)	5.1	76.8	n.a.

Note: $n = 79$, with 42 choices per person. “Fit” is the share of choices (in those subjects in whose behaviour it fits best) consistent with the rule in question. The “fit premium” is the difference between the share of these subjects’ choices explained by the given rule and the share of these subjects’ choices explained by the next-best-fitting rule.

Table 7. Subjects’ rationales.

Choice \ Rationales	Similar in terms of health gain		Wholly dissimilar			On average, subjects who use similarity heuristic switch to the following in wholly dissimilar choices	
	$A \vee B$	$C \vee D$	$A \vee E$	$G \vee F$	$Q \vee P$	Share of all subjects	Share of subjects who appeal to similarity
Similarity	41.8	44.3	2.5	0.0	0.0		
Worse off	30.4	36.7	63.3	79.7	69.6	30.2	70.1
Greater number	8.9	3.8	19.0	12.7	8.9	8.6	20.1
Total utility	6.7	3.8	3.8	1.3	12.7	0.8	2.0
No rationale	13.9	11.4	11.4	6.3	8.9	3.4	7.8
Total	100.0	100.0	100.0	100.0	100.0	43.0	100.0

Note: $n = 79$. Numbers are percentages of the whole subject population, except the final column, which lists percentages of the population that appealed to similarity to explain at least one of their choices.