

1     **A DEFENSE OF PROSPECT THEORY IN BERNHEIM &**  
2                     **SPRENGER’S EXPERIMENT (AL<sub>1</sub>)**

3     MOHAMMED ABDELLAOU<sup>a</sup>, CHEN LI<sup>b</sup>, PETER P. WAKKER<sup>b</sup>, & GEORGE WU<sup>c</sup>

4             a: HEC Paris and CNRS, Jouy-en-Josas, France, abdellaoui@hec.fr

5             b: Erasmus School of Economics, Erasmus University Rotterdam, Rotterdam, the  
6                     Netherlands, c.li@ese.eur.nl; Wakker@ese.eur.nl

7             c: Booth School of Business, University of Chicago, Chicago, USA,  
8                     wu@chicagobooth.edu

9                     20 August, 2020

10

11             Bernheim and Sprenger (2020, *Econometrica*) presented experimental evidence  
12     aimed to falsify rank dependence (and, thus, prospect theory). We argue that their  
13     experiment captured heuristics and not preferences. The same tests, but with  
14     procedures that avoid heuristics, have been done before, and they confirm rank  
15     dependence. Many other violations of rank dependence have been published before.  
16     Bernheim and Sprenger recommend rank-independent probability weighting with  
17     complexity aversion, but this is theoretically unsound and empirically invalid. In view  
18     of its many positive results, prospect theory with rank dependence remains the best  
19     model of probability weighting and the existing model that works best for applied  
20     economics.

21

22     JEL-CLASSIFICATION: D81, C91

23     KEYWORDS: prospect theory; rank dependence; complexity aversion

24

25

1

## 1. INTRODUCTION

2 Bernheim & Sprenger (2020) (BS henceforth) claim to have experimentally falsified  
 3 rank-dependent probability weighting and, hence, Tversky & Kahneman's (1992)  
 4 (cumulative) prospect theory (CPT). We dispute this claim. The main problem in BS's  
 5 experiment is that their stimuli are too complex while the stakes are too low. Many  
 6 preceding papers have argued that such stimuli lead to responses based on heuristics  
 7 rather than preferences; see §§4-5. Replications are desirable that circumvent these  
 8 issues with filler questions, visual aids, and larger stakes (outcome *differences*).  
 9 Fortunately, such replications have already been provided in the literature: Weber &  
 10 Kirsner (1997) for BS's first experiment (our Eq. 5), and Diecidue, Wakker, &  
 11 Zeelenberg (2007) for BS's second experiment (our Eq. 7). Both studies confirmed  
 12 rank dependence, showing that the findings of BS are not robust.<sup>1</sup>

13 BS suggest a theory of rank-independent weighting. However, such weighting is  
 14 unsound, more than commonly thought (§2.4). In their experimental measurement, BS  
 15 overlook that a common power of probability weighting and utility cannot be  
 16 identified from their stimuli. BS further suggest that the preference functional may  
 17 depend on the number of outcomes of a lottery. Many papers have discussed this idea  
 18 (§6.3). Tversky & Kahneman (1992 p. 317), for instance, argued that such  
 19 dependence, as well as similar effects, do not lend themselves to formal analysis.  
 20 Indeed, the idea never became popular in economics. BS argue for an aversion to  
 21 many gains, but most empirical studies find the opposite: a preference for many gains  
 22 (§6.3).

23 BS criticize a commonly-used statistical technique, but we dispute their criticism  
 24 (§6.1). Possibly based in part on that criticism, they do not cite (or cite but do not  
 25 discuss) much preceding literature. We understand that the rank-dependent stream is  
 26 too big to completely survey, and thus add several references to preceding violations  
 27 of rank dependence (§6.2, §6.4). Some are like those of BS but more serious.

28 Beyond our critique of BS's experiments, we argue that models such as prospect  
 29 theory should be evaluated by considering the body of tests of these theories, with no  
 30 single test invalidating a model. Besides the many aforementioned violations, many

---

<sup>1</sup> It is understandable that papers in field journals of over a decade ago have not been widely known.

1 more studies have supported rank-dependent probability weighting. Its imperfections  
 2 notwithstanding, prospect theory is the best model with probability weighting for use  
 3 in modeling economic phenomena such as insurance or asset pricing (Barberis 2013;  
 4 Fehr-Duda & Epper 2012). Quiggin (1982 Eq. 10) showed that, under mild  
 5 assumptions, rank dependence is the only probability weighting model that does not  
 6 violate stochastic dominance.

7 In sum, in our critique of BS, we put forth these three considerations – whether  
 8 the experiment measures preferences or heuristics; how the new empirical evidence  
 9 adds to prior empirical evidence; and the performance of the model in applied settings  
 10 – as essential for making sense of empirical evidence and the models that they test.  
 11

## 12 2. THREE PROBLEMS FOR BS’S TREATMENT OF 1979 13 PROSPECT THEORY

14 By

$$15 \quad (p: X, q: Y, 1 - p - q: Z), \quad (1)$$

16 called *lottery*, we denote a probability distribution on  $\mathbb{R}^+$  that assigns probability  $p$  to  
 17  $X$ , probability  $q$  to  $Y$ , and probability  $1 - p - q$  to  $Z$  ( $p \geq 0, q \geq 0, p + q \leq 1$ ). In  
 18 what follows, we use BS’s notation and terminology as much as possible.<sup>2</sup> BS only  
 19 consider lotteries with three or fewer outcomes, which are all gains ( $\geq 0$ ). Fewer  
 20 outcomes result if some of the probabilities in Eq. 1 are 0. By  $u: \mathbb{R}^+ \rightarrow \mathbb{R}^+$  we denote  
 21 a *utility function* (or value function). It is assumed to be strictly increasing and  
 22 continuous, with  $u(0) = 0$ . By  $\pi: [0,1] \rightarrow [0,1]$  we denote a *weighting function*. It is  
 23 assumed to be strictly increasing with  $\pi(0) = 0$  and  $\pi(1) = 1$ . For original prospect  
 24 theory (PT; Kahneman & Tversky 1979), BS propose the following evaluation of the  
 25 lottery in Eq. 1:

$$26 \quad \pi(p)u(X) + \pi(q)u(Y) + \pi(1 - p - q)u(Z). \quad (2)$$

---

<sup>2</sup> We do not use BS’s notation of lotteries because their use of braces to denote arrays rather than sets is unconventional.

1 This is an understandable proposal that has often been made in the literature but, and  
 2 this is the first problem, it is not correct strictly speaking. Eq. 2 has been known as  
 3 *separable prospect theory*. The correct formula for PT, where we may assume  
 4  $X > Y > Z$  and  $1 - p - q > 0$ , is:

$$5 \quad \pi(p)(u(X) - u(Z)) + \pi(q)(u(Y) - u(Z)) + u(Z). \quad (3)$$

6 That is, the lowest outcome  $Z$ , called riskless by Kahneman and Tversky, should not  
 7 be weighted. The Appendix explains the background of this formula.

8 The second problem of BT's rank-independent weighting is that both formulas  
 9 (Eqs. 2 and 3) are not theoretically sound. In particular, as has often been pointed out,  
 10 they violate stochastic dominance. We add here that they can lead to violations that  
 11 are absurd in magnitude. Consider the lottery yielding outcome 0 with probability  
 12 0.01, and outcome  $1 + j \times 10^{-5}$  with probability 0.01 for  $j = 1, \dots, 99$ . The certainty  
 13 equivalent of the lottery, with the parametric estimates of Tversky & Kahneman  
 14 (1992) and under the extensions of both Eq. 2 and Eq. 3, is 6.90, which exceeds the  
 15 maximal outcome of the lottery by a factor of more than 6. This does not make any  
 16 sense.<sup>3</sup> Further, Rieger & Wang (2008) showed that any extension of the original PT  
 17 to continuous distributions is problematic, depending on  $\pi$  only through  $\pi'(0)$  and  
 18 depending much on the particular discrete approximations chosen.

19 One of the main empirical findings concerns the overweighting of extreme  
 20 outcomes (Fehr-Duda & Epper 2012; l'Haridon & Vieider 2019; Luce & Suppes 1965  
 21 §4.3; Starmer 2000). It fits well with rank dependence, but cannot be accommodated  
 22 with rank-independent weighting. For all the aforementioned reasons, the rank-  
 23 independent versions of PT have been generally abandoned in favor of rank-  
 24 dependence (Barberis 2013 p. 174).

25 BS aimed to measure probability weighting and utility. To do so, they only  
 26 considered lotteries with one nonzero outcome in their first and second experiment,  
 27 probably because they used these measurements both in their rank-dependent and  
 28 rank-independent weighting analyses. However, this gives rise to the third problem: a  
 29 joint power of probability weighting and utility cannot be identified from these  
 30 stimuli. Thus,  $\pi(p)u(x)$  is empirically indistinguishable from  $\pi(p)^r u(x)^r$  for any

---

<sup>3</sup> Normalizing decision weights, as in BS (p. 1402 top) does not help (Wakker 2010 pp. 275-276).

1  $r > 0$  (Cohen & Jaffray 1988 Eq. 7a). For this reason, Fehr-Duda & Epper (2012 p.  
2 583) strongly advised against using only such stimuli. BS find that

$$3 \quad (p: x, 1 - p: 0) \rightarrow \frac{p^{0.715}}{(p^{0.715} + (1-p)^{0.715})^{1/0.715}} x^{0.941}$$

4 fits their data best in Experiment 1. But  $\left(\frac{p^{0.715}}{(p^{0.715} + (1-p)^{0.715})^{1/0.715}}\right)^{\frac{1}{0.941}} x$  fits their data  
5 equally well. In particular, the power family that they assume for utility can never rule  
6 out linear (or convex or concave) utility as best fitting. By taking  $r$  above small or  
7 large enough,  $\pi$  can be as high or low as desired, violating the inverse-S shape of the  
8 probability weighting function that BS claim as optimal and that is commonly found.

9 Similarly, in BS's second experiment,  $\left(\frac{p^{0.766}}{(p^{0.766} + (1-p)^{0.766})^{1/0.766}}\right)^{\frac{1}{0.982}} x$  fits the data  
10 equally well as their claimed optimal fit. Hence, BS could not really measure their  
11 model empirically.

12 Because of the three aforementioned problems, we will not discuss BS's analyses  
13 of PT further, including their rank-independent probability weighting. We instead  
14 focus on BS's analyses of rank dependence henceforth.  
15

### 16 3. DETERMINISTIC ANALYSIS OF BS'S EXPERIMENTS

17 This section gives some preparatory mathematical definitions. It also displays an  
18 assumption of linear utility that will be central in later discussions and is by itself  
19 reasonable. We assume CPT, with the following evaluation of lotteries:

$$20 \quad (p: X, q: Y, 1 - p - q: Z) \rightarrow w_X u(X) + w_Y u(Y) + w_Z u(Z). \quad (4)$$

21 Here,  $u$  is as above, and  $w_X$ ,  $w_Y$ , and  $w_Z$  are *decision weights*. Decision weights are  
22 *rank-dependent*. For example,  $w_X$  depends on whether  $X$  is the best, middle, or worst  
23 outcome. We follow BS in using the term rank informally and in not expressing rank  
24 dependence in notation.<sup>4</sup>

25 BS's first and main experiment concerns indifferences of the form

---

<sup>4</sup> Wakker (2010) formalized ranks and rank dependence.

$$(p: X, q: Y, 1 - p - q: Z) \sim (p: X, q: Y + m, 1 - p - q: Z - k). \quad (5)$$

1  $X$ , the common outcome, is varied across BS's main experiment, with  $m$  and  $k$  so  
 2 small that the ranking of outcomes does not change (i.e., "comonotonicity" is  
 3 satisfied). Throughout,  $Y = 24$ ,  $Z = 18$ ,  $m = 5$ , and  $q = 0.3$ . Price lists are used to  
 4 elicit  $k$  (called an *equalizing reduction* by BS) for three different values of  $p$ :  $p = 0.1$ ,  
 5  $p = 0.4$ , or  $p = 0.6$ .

6 BS used seven values of  $X$ . In some instances,  $X$  is the best ranked outcome  
 7 ( $X = 34$ ,  $X = 32$ , or  $X = 30$ ); in other cases,  $X$ , which we denote  $X'$ , is ranked in the  
 8 middle ( $X' = 23$ ,  $X' = 21$ , or  $X' = 19$ ). When  $X$  is ranked best, the weights for  
 9 outcomes  $Y$  and  $Y + m$  are denoted  $w_Y$ . When  $X$  is ranked middle, they are denoted  
 10  $w_{Y'}$ .

11 Under linear utility,  $\frac{w_Y}{w_Z} = \frac{k}{5}$  and  $\frac{w_{Y'}}{w_Z} = \frac{k'}{5}$ . The ratio

$$\frac{k'}{k} = \frac{w_{Y'}}{w_Y} \quad (6)$$

12 captures the proportional change of the decision weight and, hence, rank dependence.  
 13 This ratio, or its log, is used in BS's analysis.

14 BS repeatedly claim that they can have Eq. 6 for all differentiable utility  
 15 functions. However, this claim is based on marginal rates of substitution involving  
 16 infinitesimal changes  $m$  and  $k$ , which cannot be implemented empirically.<sup>5</sup>  
 17 Empirically, we have to work with the following, reasonable, assumption:

18 ASSUMPTION 1 [linear utility for moderate outcome changes]. For outcome changes  
 19 within a small interval  $[A, B]$ , utility is approximately linear.  $\square$

20 More precisely, for Eq. 6, it can be seen that linearity of utility is used on the interval  
 21  $[\min\{18 - k, 18 - k'\}, 18]$ . Assumption 1 provides a good approximation for all

---

<sup>5</sup> BS's Footnote 13 even claims validity for infinitesimals for every strictly increasing continuous utility, dispensing with differentiability. However, this is not correct. For singular Cantor-type functions the *positive* right derivatives claimed by BS may exist *nowhere*, let be at the points where needed. See Paradís, Viader, & Bibiloni (2001; their Theorem 3.1 and its proof are also valid for right and left derivatives).

1 common nonlinear utility functions. Empirical evidence supporting it is also  
2 abundant.<sup>6</sup>

3 BS's second experiment concerned indifferences of the form

$$4 \quad (p: X, q: Y, 1 - p - q: Z) \sim (p: X + m, q: Y - k, 1 - p - q: Z - k). \quad (7)$$

5 In this experiment,  $p = 0.4$  and  $q = 0.3$ , or  $p = 0.6$  and  $q = 0.2$ , with  $Y = 36$ ,  
6  $Z = 18$ , and  $m = 4$  throughout. Finally,  $X = 2, 3, 4, 20, 21, 22, 38, 39$ , or  $40$ , with  
7 price lists again used to elicit  $k$ , which BS again call an *equalizing reduction*. Under  
8 linear utility and Eq. 4,

$$9 \quad w_X = \frac{k}{m+k}. \quad (8)$$

10 Consider  $X = 4$  (with  $k$ ) and  $X' = 20$  (with  $k'$ ). In the lottery with  $X = 4$ ,  $X$  has the  
11 worst rank, whereas  $X' = 20$  has the middle rank in the corresponding lottery. The  
12 rank of  $Z$  changed from middle to worst in these two lotteries. Such rank changes  
13 affect the decision weight in Eq. 8. BS again capture the change by the ratio  $\frac{k'}{k}$ .<sup>7</sup> Eq. 8  
14 again uses Assumption 1. More precisely, it uses linearity of utility on the intervals  
15  $[\min\{18 - k, 18 - k'\}, 18]$  and  $[\min\{36 - k, 36 - k'\}, 36]$ . As argued before and  
16 supported by numerical analyses by BS, this is a reasonable assumption.

17

## 18 4. SMALL PAYOFF CHANGES: RAMSEY'S TRIFLE 19 PROBLEM AND A STATISTICAL PROBLEM

20 Ramsey (1931) pointed out a difficulty that applies to BS's implementation of  
21 Assumption 1, which we call *Ramsey's trifle problem*:

---

<sup>6</sup> See Birnbaum (2008 p. 469), Epper, Fehr-Duda, & Bruhin (2011), Homonoff (2018 p. 182),  
Kahneman & Lovallo (1993), Lopes & Oden (1999 footnote 1), Luce (2000 p. 86), Marshall (1890  
Book III), Pigou (1920 p. 785), Rabin (2000), Savage (1971 p. 786).

<sup>7</sup> More precisely, the ratio of decision weights is  $\frac{k'}{k} \times \frac{m+k}{m+k'}$  which is, roughly, a monotonic nonlinear  
transformation of  $k'/k$ . Importantly, it does not affect being larger or smaller than 1. We conjecture  
that BS still used  $k'/k$ , or its log, as index in their analysis for this reason.

1           Since it is universally agreed that money has a diminishing marginal  
 2           utility, if money bets [to measure decision weights (subjective  
 3           probabilities) through ratios] are to be used, it is evident that they should  
 4           be for as small stakes as possible. *But then again the measurement is*  
 5           *spoiled by introducing the new factor of reluctance to bother about trifles.*  
 6           [Italics added] [p. 176]

7           This trifle problem was also pointed out by Samuelson (1959 Footnote 5). In order to  
 8           approximate infinitesimal changes with perfect linearity, BS took payoff changes  
 9            $m, k$  that are very small. But these changes became too small to motivate subjects.

10           In BS's main (first) experiment, subjects completed 28 price lists, with 21 of  
 11           those (seven values of  $X$  and three different sets of probabilities) constituting the  
 12           elicitation of equalizing reductions,  $k$  or  $k'$ . Altogether, subjects answered 980  
 13           ( $21 \times 38 + 7 \times 26$ ) questions, most of which involved nearly-identical lotteries.  
 14           Although subjects earned \$26.87 on average, the effective incentive was the  
 15           possibility of getting \$4 extra with probability 0.3 at the cost of a chance at losing a  
 16           few dollars. Similar numbers and stakes appeared in BS's second experiment. It is  
 17           inconceivable that subjects, even if only subconsciously in an as-if sense, would do  
 18           anything near determining preference values of these complex lotteries, several  
 19           hundreds of times in a row, for such small stakes. Instead, it is likely that subjects,  
 20           when facing the figures of the stimuli (BS Online Appendix) will quickly recognize  
 21           the structure of Eq. 5 or 7, develop a simple algebraic heuristic and use that repeatedly  
 22           (combined with the usual noise) to quickly get through the experiment (§5). Smith  
 23           (1982) posited a dominance requirement for experimental economics: the rewards  
 24           should dominate subjective costs. Wilcox (1993) confirmed empirically that good  
 25           incentives are necessary for complex stimuli. These requirements are violated by BS.  
 26           von Winterfeldt & Edwards (1982) reviewed a number of studies that showed that  
 27           subjects use simple strategies in situations with inadequate incentives.

28           We also note a statistical problem concerning preferences (as opposed to BS's  
 29           data). BS's §2.1 reports a deterministic CPT preference analysis of their stimuli that  
 30           uses ratios  $\frac{k}{k'}$ , where  $k$  and  $k'$  are small in an absolute sense relative to the other  
 31           numbers in the stimuli. We note here that  $18 - k$  and  $18 - k'$  were the values actually  
 32           elicited. Small relative errors in these give large relative errors in  $k, k'$ . Hence, ratios  
 33            $\frac{k}{k'}$  are very vulnerable to noise. As BS emphasize throughout, it is important to reckon  
 34           with noise beyond a deterministic analysis. It would have been of special interest to



1 analyze the role of plausible noise *in preferences* in their §2.1. Adding an error term  
 2 to their CPT values (as in BS’s Eq. 5 and Footnote 60) affects the certainty  
 3 equivalents of the overall lotteries by some dollars. Given the complexity of their  
 4 three-outcome lotteries and Ramsey’s trifle problem, such errors in preference values  
 5 are plausible. This leads to errors in the measured  $k, k'$  that may readily make them  
 6 approximate 0 (no negative answers were possible). If such errors occur with  
 7 probabilities exceeding 0.05, then the confidence intervals of the ratios  $\frac{k}{k'}$  span the  
 8 whole  $\mathbb{R}^+$ . Then BS’s analysis may lack the statistical power to reject any hypothesis  
 9 *about preference*, be it rank dependence or rank independence.

10 The above statistical analysis was back-of-the-envelope, for illustrative purposes.  
 11 It shows that a fully elaborated power analysis, based on adding plausible error  
 12 models *for preferences* to the calibrated CPT models used throughout BS’s paper,  
 13 would have provided useful insights. It would have shown if the claim in their abstract  
 14 “Conventional calibrations of CPT preferences imply that the percentage change in  
 15 probability weights should be an order of magnitude larger than we observe,” and  
 16 claimed nonoverlapping confidence intervals (BS p. 1366 middle; p. 1382 l. 12; p.  
 17 1388), can hold statistically for noisy-preference calibration models. It would also  
 18 show whether the variances found in their data may at all represent preferences rather  
 19 than heuristics, as we argue. In general, power analyses are best done prior to  
 20 observing data. For brevity, we do not elaborate on them.

21

## 22 5. HEURISTICS IN BS’S EXPERIMENTS

23 The data that BS found did not exhibit the volatility suggested at the end of the  
 24 preceding section. BS obtained stable patterns giving statistical power and tight  
 25 confidence intervals. Our claim is that this is because their experiment did not  
 26 measure preferences, and did not, even in an as-if sense, speak to Eqs. 2, 3, or 4.  
 27 Instead, subjects faced with hundreds of choices of complex and nearly-identical  
 28 lotteries, for a one-time trifle reward received with some probability, develop simple  
 29 algebraic heuristics to get through the experiment.

30 Many studies have shown that multiple repetitions of complex tasks can lead to  
 31 stable but invalid patterns, in our case heuristics instead of preferences. Ariely,

1 Loewenstein, & Prelec (2001) call it coherent arbitrariness, while Loomes, Starmer, &  
 2 Sugden (2003) call it the shaping hypothesis. See also Baron et al. (2001 p. 3 l. -2),  
 3 Carlin (1992 p. 219), Dolan & Stalmeier (2003), and Hardisty et al. (2013).

4 This main heuristic in BS's first and main experiment, based on Eq. 5, was  
 5 cancellation: subjects ignored the common outcome  $X$ , not because of preference, but  
 6 only as a heuristic to simplify the experimental task. It precludes rank dependence. BS  
 7 (p. 1366) acknowledge this problem, citing Wu (1994) and Weber & Kirsner (1997).  
 8 Cancellation has been widely documented in the literature (see below). Weber &  
 9 Kirsner (1997 top of p. 57) showed that Wakker, Erev, & Weber (1994) suffered from  
 10 cancellation, explaining the absence of rank dependence there. All these papers used  
 11 stimuli as in BS's first experiment (Eq. 5). In a treatment using such stimuli while  
 12 avoiding cancellation (by asking for pricing rather than direct choice) and other  
 13 heuristics, Weber & Kirsner did find rank dependence. The problem of cancellation in  
 14 Wakker, Erev, & Weber's design reappears in BS's first experiment.

15 A general finding is that, the more saliently the common outcome is displayed,  
 16 the stronger cancellation is.<sup>8</sup> In the stimuli of BS's first experiment, throughout and  
 17 invariably, the most left column of lotteries displays the common outcome, each time  
 18 spanning the whole page (BS Online Appendix). This is as salient as it can be.  
 19 Heuristics were further facilitated by the following features of the stimuli, that can be  
 20 inferred from the experimental instructions provided online. Invariably, the middle  
 21 column displayed the outcomes 24 (for the left lottery) and 29 (for the right lottery),  
 22 always with the same probability 0.3. Further, the right column always displayed  
 23 outcomes \$18 (for the left lottery) and \$18 -  $k$  (for the right lottery), with the same  
 24 probability vector. Probabilities of those outcome always decreased in the same order  
 25 in all blocks of three.

26 To avoid cancellation, BS (§5.3) carried out a second, smaller experiment, based  
 27 on Eq. 7. Now there is no common outcome to be cancelled. They also had a bigger  
 28 variation in outcomes  $X$ , which has two advantages. First, it makes it harder for  
 29 subjects to develop some simple heuristics that ignore nonvarying outcomes. Second,

---

<sup>8</sup> We give one reference from every decade: Kahneman & Tversky (1979 p. 274—on their Figures 1 versus 2); Keller (1985); Kashima & Maher (1995); Birnbaum (2008 p.481 ff.) ; Schneider, Leland, & Wilcox (2018); Blavatsky, Ortmann, & Panchenko (2020 —compound vs. lottery).

1 middle ranks are not only changed into best ranks, where rank dependence is known  
 2 to be weak (§6.2), but also into worst ranks, where rank dependence is known to be  
 3 stronger (Fehr-Duda & Epper 2012 end of §2; Wakker 2010 p. 227; Weber & Kirsner  
 4 1997 p. 58).

5 The format in BS's second experiment has not been used much before and, hence,  
 6 less is known about presence or absence of heuristics. Despite the aforementioned  
 7 advantages, we still think that heuristics were measured and not preferences, because  
 8 many problems of BS's Experiment 1 remain in their Experiment 2. The complexity  
 9 of the lotteries has worsened due to the absence of a common outcome. This augments  
 10 Ramsey's trifle problem. The statistical problem of ratios of small numbers also  
 11 remains. The layout of the stimuli (their Online Appendix, Figure 5) facilitates the  
 12 heuristics with, again, the same format for hundreds of choices over several pages,  
 13 again inducing coherent arbitrariness.

14 It is, for instance, heroic to think that, for the many complex lotteries and trifle  
 15 rewards, subjects would incorporate the separate values  $Y = 24, Y = 24 - k, Z =$   
 16  $18, Z = 18 - k, q,$  and  $1 - p - q$  from Eq. 7 into their valuation, even in an as-if  
 17 sense. Instead, if the values were subsumed into one concept, "lose  $k$  otherwise,"<sup>9</sup>  
 18 there is no perception of the ranking of outcomes and no scope for rank dependence.

19 Diecidue, Wakker, & Zeelenberg (2007; DWZ henceforth) used the same  
 20 equalizing reductions as in BS's second experiment, i.e., indifferences as in our Eq. 7;  
 21 see their Eq. 3.2. They thus avoided the common outcomes in Wakker, Erev, &  
 22 Weber (1994) that had been criticized by Weber & Kirsner (1997). They also used  
 23 Assumption 1 to obtain nonparametric estimations of decision weights, and their  
 24 primary purpose was also to test rank dependence (main hypothesis on p. 185; p. 195  
 25 2nd para). Their amounts  $m, k$  were considerably larger than in BS ( $m$  never below  
 26 Dfl. 20, which in present value would be about \$20). We discuss this difference. The  
 27 exact validity for general utility in BS's theorems only concerns infinitesimal  $m, k$ . In  
 28 experiments, as in BS and DWZ, one needs to take small discrete  $m, k$ , small enough  
 29 to have good approximations. However, if  $m, k$  are too small, as in BS, then one runs  
 30 into Ramsey's trifle problem, and the statistical problem at the end of §4. Hence,  $m, k$

---

<sup>9</sup> Given probability  $p$  for  $X$  and  $X + 4$ , loss  $k$  automatically occurs with complementary probability  $1 - p$ .

1 have to be larger (but not too much), as in DWZ. And one needs Assumption 1, which  
2 is still reasonable for the amounts considered (referenced above). This is the only way  
3 in which BS's approach with equalizing reductions can be used, and this is what DWZ  
4 did.

5 In the beginning of §3, DWZ explain that three-outcome prospects are too  
6 difficult to evaluate in general; see also DWZ (p. 181, 3rd & 4th para). Hence, they  
7 used a visual design (their Figure 1) to facilitate these choices,<sup>10</sup> developed in  
8 extensive pilot studies with debriefings to identify and then avoid the major heuristics  
9 used by subjects (see their p. 188 last two paras; p. 194 last para; p. 195 last para).  
10 They used filler questions, and more variations in outcomes, to further reduce  
11 heuristics. Following Weber & Kirsner (1997), DWZ minimized all aforementioned  
12 problems in the design of BS.<sup>11</sup> Podsakoff, MacKenzie, & Podsakoff (2012) survey  
13 general techniques for reducing measurement-instrument biases. DWZ emphasized  
14 that avoiding heuristics is desirable to increase statistical power (p. 188 last line; p.  
15 195 3rd para). Their findings supported rank dependence. That is, there were  
16 significant changes of decision weight if ranks of events changed (p. 192 last para), in  
17 agreement with other findings of rank dependence in the literature.

18 DWZ also found violations of comonotonic independence, that is, of rank  
19 dependence. Decision weights sometimes changed even though ranks did not. BS's  
20 finding is reversed in the sense that their decision weights did not change even if  
21 ranks did. Whereas BS's finding can be taken as a special case of CPT (expected  
22 utility), DWZ's finding is more serious because it cannot be reconciled with any  
23 version of CPT. Wu (1994) showed large violations of rank-dependence by testing

---

<sup>10</sup> This is commonly done for lotteries with three or more outcomes, by Weber & Kirsner (1997) and most others. Lola Lopes, specialized in multi-outcome lotteries, developed special visual designs (Lopes & Oden 1999; Fennema & Wakker 1997).

<sup>11</sup> Two further differences between DWZ and BS are as follows. First, DWZ used events with unknown probabilities rather than probabilities. This does not affect the theoretical working of rank dependence. See DWZ (p. 185 ll. 7-8) and Wakker (2010 Figure 7.4.1 versus 10.4.1). Second, an improvement of one prospect was not compensated by worsening that same prospect elsewhere, but instead by improving the other prospect. This avoids confounds due to differences between improvements and worsenings.

1 comonotonic independence. In particular, he showed that heuristics such as  
2 cancellation can sometimes generate patterns inconsistent with prospect theory.

3 It could be argued that the different findings of DWZ regarding rank dependence  
4 are because they considered uncertainty rather than risk. That deviations from  
5 expected utility are rank dependent for uncertainty but not for risk. However, we find  
6 this explanation implausible. We think that the differences are because DWZ's design,  
7 unlike BS's design, reduced heuristics and errors so that they achieved statistical  
8 power to test *preference* conditions.

9

## 10 6. FURTHER COMMENTS

### 11 6.1. BS'S CRITICISM OF COMMON STATISTICAL TESTS

12 BS (§2.3) criticize tests that compare numbers of violations of predictions  
13 (preference axioms), which are commonly used throughout decision theory and in  
14 other empirical sciences. BS erroneously claim that such tests would not assume any  
15 ("parametric") noise model, writing:

16 These types of frequency comparisons raise two difficulties, both  
17 stemming from the fact that the results are difficult to interpret without a  
18 parametric model of noisy choice. First, the premise of the approach—that  
19 violation frequencies are *necessarily* higher for invalid axioms—is flawed.  
20 [italics added] [p. 1376]

21 See also BS (p. 1367 l. -3 and p. 1376 2nd para). However, the common tests are  
22 statistical, and statistical tests always assume a noise model, contrary to BS's  
23 "necessarily" claim. These tests only assume that the aforementioned higher  
24 frequencies are likely, not necessary. Following up, BS claim the following two  
25 difficulties.

26 The first claimed difficulty concerns a counterexample where all choices testing  
27 one preference axiom are close to indifference and subject to much noise, whereas  
28 those for the other preference axiom are all far from indifference, with little noise. BS  
29 claim that this would invalidate the common tests. We argue that one cannot disregard  
30 a whole stream of literature based on one artificial counter-example. For every data  
31 set and statistical analysis based on a noise model, one can specify an alternative noise  
32 model that invalidates the analysis. However, to serve as a good counterexample, the  
33 alternative noise model should be plausible. Every good test based on counting

1 violations was done in a design where BS's first difficulty was not plausible. Note  
 2 also that there have been many such tests in the literature. Even if one by accident was  
 3 as in BS's example, then this does not invalidate the whole literature based on it. We  
 4 finally point out that, even under common designs with plausible error theories,  
 5 unlikely data may arise as in BS's example, or other kinds of unlikely data, leading to  
 6 errors of type I or II. Statistical tests never claim that such errors are "necessarily"  
 7 absent; only, that they are unlikely. P-values, powers, Bayes factors, and so on capture  
 8 such unlikely exceptions.

9 BS's second claimed difficulty concerns the example at the end of their Online  
 10 Appendix B. In this example, they assume rank dependence that "inadvertently" is too  
 11 weak to affect optimal choice between any of the stimuli considered, and trembling-  
 12 hand errors depending only on preference and not on utility differences. Then CPT  
 13 and expected utility have identical predictions in this design. Tests based on numbers  
 14 of violations then indeed have no power to distinguish. But then, neither does any  
 15 other test. This trivial example cannot serve as a criticism of any test. We were unable  
 16 to understand BS's description of this example in their main text (p. 1376): "For any  
 17 given degree of rank dependence, one can construct simple examples (with constant  
 18 "distance to indifference") in which the differential between violation frequencies  
 19 falls anywhere between zero and unity." BS argue that they have refuted some widely  
 20 used statistical tests and a whole stream of literature based on it. However, their  
 21 criticisms are unfounded.

22

## 23 6.2. WEAKNESS OF RANK DEPENDENCE IN LONGSHOT EFFECT

24 In their main experiment, BS consider changes in decision weights only when the  
 25 rank changes from middle to best. It is well-known that rank dependence is not strong  
 26 there (DWZ p. 185 ll. 4-6; DWZ p. 197 l. 7). The prevailing finding is that the  
 27 weights then increase, consistent with inverse-S probability weighting and the  
 28 longshot effect. Stronger rank dependence occurs when ranks change from middle to  
 29 worst, consistent with the certainty effect. As for the change in rank considered by  
 30 BS, quite some studies found that the increase of decision weight is weak or absent.  
 31 Even, several studies found the opposite effect to be prevailing, of decreasing decision  
 32 weight, consistent with pessimistic probability weighting. See van de Kuilen &  
 33 Wakker (2011). Their Footnotes 7 & 8 survey the many other papers that found this

1 opposite effect. In view of this literature, finding no effect of rank dependence in BS's  
2 first experiment is no surprise anyhow.

3

### 4 6.3. COMPLEXITY SEEKING INSTEAD OF AVERSION

5 BS favor adding a component to risk theory that reckons with the number of  
6 outcomes of a lottery, to capture complexity aversion. Many authors have investigated  
7 proposals of this sort. Neilson (1992) proposed a formal model, but it was tested  
8 unsuccessfully by Humphrey (2001). Related models received some attention in  
9 psychology (Birnbaum 2008 p. 481; Krantz et al. 1971 Ch. 8; Luce 2000). Tversky &  
10 Kahneman (1992 p. 317) were pessimistic about modeling this phenomenon:

11 Despite its greater generality, the cumulative functional is unlikely to be  
12 accurate in detail. We suspect that decision weights may be sensitive to the  
13 formulation of the prospects, as well as to the *number*, the spacing and the  
14 level of outcomes. ... The present theory can be generalized to  
15 accommodate such effects but it is questionable whether the gain in  
16 descriptive validity, achieved by giving up the separability of values and  
17 weights, would justify the loss of predictive power and the cost of  
18 increased complexity. ... The heuristics of choice do not readily lend  
19 themselves to formal analysis because their application depends on the  
20 formulation of the problem, the method of elicitation, and the context of  
21 choice. [italics added]

22 We agree with this pessimism, and this old model never caught on in economics to  
23 our best knowledge. The volatility of empirical findings adds to our pessimism. In a  
24 literature review (Online Appendix), we found two additional empirical studies  
25 confirming complexity aversion, but seven studies finding the opposite, complexity  
26 seeking. Thus, the prevailing empirical finding is opposite to BS's model.

27

### 28 6.4. FURTHER PRECEDING FALSIFICATIONS OF PROSPECT THEORY

29 Given that prospect theory is the most tested risk theory, besides being most  
30 confirmed, it is also most violated. The keyword "PT falsified" in Wakker (2020)  
31 gives 49 papers.

32 We call special attention to a finding of Birnbaum & McIntosh (1996), which is  
33 not cited by BS. It was confirmed in several follow-up studies by Birnbaum and  
34 colleagues, surveyed by Birnbaum (2008 pp. 484-487), and found independently by  
35 Humphrey & Verschoor (2004). This finding concerns lotteries of the same format as  
36 in Eq. 5, i.e., as in BS's first experiment, with the common outcome  $X$  moved to test  
37 rank dependence. Prospect theory predicts that weights increase if ranks change from

1 middle to best or worst. BS quantitatively find no change in decision weights. The  
 2 aforementioned studies attempted to avoid heuristics and found in fact a stronger  
 3 deviation: a decrease in weight.

4 The strongest counterexample to rank dependence that we are aware of is  
 5 Machina's (2009) reflection example, confirmed empirically by l'Haridon & Placido  
 6 (2010). Ending on a positive note, Fehr-Duda & Epper (2012) and Barberis (2013)  
 7 provide surveys with many useful applications of rank dependence. See also the  
 8 impressive data sets of l'Haridon & Vieider (2019) and Ruggeri et al. (2020), and  
 9 DellaVigna's (2018) review in the context of structural behavioral economics. A  
 10 psychological justification for rank-dependent decision weights is based on cognitive  
 11 attention: individuals tend to attend to extreme outcomes (Weber & Kirsner 1997;  
 12 Pachur et al. 2018).

13

14

## 7. CONCLUSION

15 We have offered a critique of BS's experiments and analyzes. While we applaud  
 16 BS's call for more investigations of rank dependence, we attribute their observed  
 17 stability of equalizing reductions to subjects' use of heuristics, not a failure of rank  
 18 dependence as claimed by BS. Every study should be interpreted relative to a large set  
 19 of admittedly mixed empirical results. Prospect theory is an imperfect theory, as will  
 20 be every theory that aims to make sense of the complex problem of how people make  
 21 decisions. Nevertheless, for now, we see prospect theory as the most tractable and  
 22 best performing model with nonlinear probability weighting.

23

## 24 APPENDIX A: BACKGROUND OF PROSPECT THEORY

25

### FORMULA OF 1979

26 Kahneman & Tversky (1979) only wrote Eq. 3 explicitly for at most two nonzero  
 27 outcomes, i.e., when either  $Z = 0$  (their Eq. 1) or  $p = 0$  (their Eq. 2). This has led to  
 28 much confusion about Eqs. 2 and 3 in the literature. We quote their verbal description  
 29 to show that Eq. 3 is correct.



1 ... prospects are segregated into two components: (i) the riskless  
 2 component, i.e., the minimum gain ... which is certain to be obtained or  
 3 paid; (ii) the risky component, i.e., the additional gain[s] ... which is[are]  
 4 actually at stake.... That is, the value of a strictly positive ... prospect  
 5 equals the value of the riskless component plus the value-difference  
 6 between the outcomes, multiplied by the weight associated with the more  
 7 extreme outcome[s]. The essential feature ... is that a decision weight is  
 8 applied to the value difference ... which represents the risky component of  
 9 the prospect, but not to ... the riskless component. (Kahneman & Tversky.  
 10 1979 p. 276).

11 We added the texts between square brackets to extend to plurality, with multiple non-  
 12 minimal outcomes.

13 From the case  $p = 0$  it appears that Eq. 2 above, claimed by BS, cannot be  
 14 correct, as the riskless outcome  $Z$  should not be weighted. Eq. 3 is the only natural  
 15 extension of the formulas given by Kahneman & Tversky. That Eq. 3 is correct also  
 16 appears from Kahneman & Tversky (1975 p. 18), a first working paper version of  
 17 their 1979 paper. They do explicitly write the formula of PT for multiple outcomes  
 18 there. They treat the value function somewhat differently than in their 1979 paper,  
 19 taking utilities of differences rather than differences of utilities.<sup>12</sup> But they treat  
 20 probability weighting as in Eq. 3. They also emphasize (their p. 12) that the riskless  
 21 outcome should be as in our Eq. 3 (not weighted) rather than in our Eq. 2 (where it is  
 22 weighted). Whenever relevant, they pointed out that the riskless outcome should not  
 23 be weighted (Tversky & Kahneman 1981 Footnote 5; Tversky & Kahneman 1986  
 24 Footnote 2).

25 BS (footnote 11) cite Camerer & Ho (1994) for Eq. 2. However, Camerer & Ho  
 26 (1994) used a different term, separable prospect theory, for Eq. 2. Their endnote 16  
 27 pointed out that it deviates from prospect theory for strictly positive lotteries. BS  
 28 (footnote 11) also cite Fennema & Wakker (1997) for Eq. 2. However, Fennema &  
 29 Wakker (p. 54) pointed out that this equation should be used for mixed prospects,  
 30 which assign positive probabilities to both gains and losses. Those are among what  
 31 Kahneman & Tversky called regular prospects, and then the analog of Eq. 2 is indeed  
 32 correct (Kahneman & Tversky 1975 p. 18).

33

---

<sup>12</sup> The latter is preferable because it can be applied to nonquantitative outcomes.

## REFERENCES

- 1
- 2 Ariely, Dan, George F. Loewenstein, & Drazen Prelec (2001) “ ‘Coherent  
3 Arbitrariness’: Stable Demand Curves without Stable Preferences,” *Quarterly*  
4 *Journal of Economics* 118, 73–106.
- 5 Barberis, Nicholas (2013) “Thirty Years of Prospect Theory in Economics: A Review  
6 and Assessment,” *Journal of Economic Perspectives* 27, 173–195.
- 7 Baron, Jonathan, Zhijun Wu, Dallas J. Brennan, Christine Weeks, & Peter A. Ubel  
8 (2001) “Analog Scale, Magnitude Estimation, and Person Trade-Off as Measures  
9 of Health Utility: Biases and Their Correction,” *Journal of Behavioral Decision*  
10 *Making* 14, 17–34.
- 11 Bernheim, B. Douglas & Charles Sprenger (2020) “On the Empirical Validity of  
12 Cumulative Prospect Theory: Experimental Evidence of Rank-Independent  
13 Probability Weighting,” *Econometrica* 88, 1363–1409.
- 14 Birnbaum, Michael H. (2008) “New Paradoxes of Risky Decision Making,”  
15 *Psychological Review* 115, 463–501.
- 16 Birnbaum, Michael H. & William R. McIntosh (1996) “Violations of Branch  
17 Independence in Choices between Gambles,” *Organizational Behavior and*  
18 *Human Decision Processes* 67, 91–110.
- 19 Blavatsky, Pavlo, Andreas Ortmann, & Valentyn Panchenko (2020) “On the  
20 Experimental Robustness of the Allais Paradox.” *American Economic Journal:*  
21 *Microeconomics*, forthcoming.
- 22 Camerer, Colin F. & Teck-Hua Ho (1994) “Violations of the Betweenness Axiom and  
23 Nonlinearity in Probability,” *Journal of Risk and Uncertainty* 8, 167–196.
- 24 Carlin, Paul S. (1992) “Violations of the Reduction and Independence Axioms in  
25 Allais-Type and Common-Ratio Experiments,” *Journal of Economic Behavior*  
26 *and Organization* 19, 213–235.
- 27 Cohen, Michèle & Jean-Yves Jaffray (1988) “Certainty Effect versus Probability  
28 Distortion: An Experimental Analysis of Decision Making under Risk,” *Journal*  
29 *of Experimental Psychology: Human Perception and Performance* 14, 554–560.
- 30 DellaVigna, Stefano (2018) “Structural Behavioral Economics.” In B. Douglas  
31 Bernheim, Stefano DellaVigna, & David Laibson (eds.) *Handbook of Behavioral*  
32 *Economics; Volume 2*, 613–723, Elsevier, Amsterdam.

- 1 Diecidue, Enrico, Peter P. Wakker, & Marcel Zeelenberg (2007) “Eliciting Decision  
2 Weights by Adapting de Finetti’s Betting-Odds Method to Prospect Theory,”  
3 *Journal of Risk and Uncertainty* 34, 179–199.
- 4 Dolan, Paul & Peep Stalmeier (2003) “The Validity of Time Trade-Off Values in  
5 Calculating QALYs: Constant Proportional Time Trade-Off versus the  
6 Proportional Heuristic,” *Journal of Health Economics* 22, 445–458.
- 7 Epper, Thomas, Helga Fehr-Duda, & Adrian Bruhin (2011) “Viewing the Future  
8 through a Warped Lens: Why Uncertainty Generates Hyperbolic Discounting,”  
9 *Journal of Risk and Uncertainty* 43, 163–203.
- 10 Fehr-Duda, Helga & Thomas Epper (2012) “Probability and Risk: Foundations and  
11 Economic Implications of Probability-Dependent Risk Preferences,” *Annual  
12 Review of Economics* 4, 567–593.
- 13 Fennema, Hein & Peter P. Wakker (1997) “Original and Cumulative Prospect Theory:  
14 A Discussion of Empirical Differences,” *Journal of Behavioral Decision Making*  
15 10, 53–64.
- 16 Hardisty, David J., Katherine F. Thompson, David H. Krantz, & Elke U. Weber  
17 (2013) “How to Measure Time Preferences: An Experimental Comparison of  
18 Three Methods,” *Judgment and Decision Making* 8, 214–235.
- 19 Homonoff, Tatiana A. (2018) “Can Small Incentives Have Large Effects? The Impact  
20 of Taxes versus Bonuses on Disposable Bag Use,” *American Economic Journal:  
21 Economic Policy* 10, 177–210.
- 22 Humphrey, Steven J. (2001) “Are Event-Splitting Effects Actually Boundary  
23 Effects?,” *Journal of Risk and Uncertainty* 22, 79–93.
- 24 Humphrey, Steven J. & Arjan Verschoor (2004) “Decision-Making under Risk among  
25 Small Farmers in East Uganda,” *Journal of African Economies* 13, 44–101.
- 26 Kahneman, Daniel & Dan Lovallo (1993) “Timid Choices and Bold Forecasts: A  
27 Cognitive Perspective on Risk Taking,” *Management Science* 39, 17–31.
- 28 Kahneman, Daniel & Amos Tversky (1975) “Value Theory: An Analysis of Choices  
29 under Risk,” paper presented at the ISRACON conference on Public Economics,  
30 Jerusalem, 1975.
- 31 Kahneman, Daniel & Amos Tversky (1979) “Prospect Theory: An Analysis of  
32 Decision under Risk,” *Econometrica* 47, 263–291.
- 33 Kashima, Yoshihisa & Patrick Maher (1995) “Framing of Decisions under  
34 Ambiguity,” *Journal of Behavioral Decision Making* 8, 33–49.

- 1 Keller, L. Robin (1985) “An Empirical Investigation of Relative Risk Aversion,”  
 2 *IEEE Transactions on systems, Man, and Cybernetics*, SMC-15, 475–482.
- 3 Krantz, David H., R. Duncan Luce, Patrick Suppes, & Amos Tversky (1971),  
 4 “*Foundations of Measurement, Vol. I (Additive and Polynomial*  
 5 *Representations)*.” Academic Press, New York. (2<sup>nd</sup> edn. 2007, Dover  
 6 Publications, New York.)
- 7 l’Haridon, Olivier & Laetitia Placido (2010) “Betting on Machina's Reflection  
 8 Example: An Experiment on Ambiguity,” *Theory and Decision* 69, 375–393.
- 9 l’Haridon, Olivier & Ferdinand Viieder (2019) “All over the Map: A Worldwide  
 10 Comparison of Risk Preferences,” *Quantitative Economics* 10, 185–215.
- 11 Loomes, Graham, Chris Starmer, & Robert Sugden (2003) “Do Anomalies Disappear  
 12 in Repeated Markets?,” *Economic Journal* 113, C153–C166.
- 13 Lopes, Lola L. & Gregg C. Oden (1999) “The Role of Aspiration Level in Risky  
 14 Choice: A Comparison of Cumulative Prospect Theory and SP/A Theory,”  
 15 *Journal of Mathematical Psychology* 43, 286–313.
- 16 Luce, R. Duncan (2000) “*Utility of Gains and Losses: Measurement-Theoretical and*  
 17 *Experimental Approaches*.” Lawrence Erlbaum Publishers, London.
- 18 Luce, R. Duncan & Patrick Suppes (1965) “Preference, Utility, and Subjective  
 19 Probability.” In R. Duncan Luce, Robert R. Bush, & Eugene Galanter (eds.)  
 20 *Handbook of Mathematical Psychology*, Vol. III, 249–410, Wiley, New York.
- 21 Machina, Mark J. (2009) “Risk, Ambiguity, and the Rank-Dependence Axioms,”  
 22 *American Economic Review* 99, 385–392.
- 23 Marshall, Alfred (1890) “*Principles of Economics*.” 8<sup>th</sup> edn. 1920 (9<sup>th</sup> edn. 1961),  
 24 MacMillan, New York.
- 25 Neilson, William S. (1992) “Some Mixed Results on Boundary Effects,” *Economics*  
 26 *Letters* 39, 275–278.
- 27 Pachur, Thorsten, Michael Schulte-Mecklenbeck, Ryan O. Murphy, & Ralph Hertwig  
 28 (2018) “Prospect Theory Reflects Selective Allocation of Attention,” *Journal of*  
 29 *Experimental Psychology: General* 147, 147–169.
- 30 Paradís, Jaumnine, Pelegri Viader, & Lluís Bibiloni (2001) “The Derivative of  
 31 Minkowski's  $\varphi(x)$  Function,” *Journal of Mathematical Analysis and Applications*  
 32 253, 107–125.
- 33 Pigou, Arthur C. (1920) “*The Economics of Welfare*.” (edn. 1952: MacMillan,  
 34 London.)

- 1 Podsakoff, Philip M., Scott B. MacKenzie, & Nathan P. Podsakoff (2012) “Sources of  
2 Method Bias in Social Science Research and Recommendations on How to  
3 Control It,” *Annual Review of Psychology* 63, 539–69.
- 4 Quiggin, John (1982) “A Theory of Anticipated Utility,” *Journal of Economic  
5 Behaviour and Organization* 3, 323–343.
- 6 Rabin, Matthew (2000) “Risk Aversion and Expected-Utility Theory: A Calibration  
7 Theorem,” *Econometrica* 68, 1281–1292.
- 8 Ramsey, Frank P. (1931) “Truth and Probability.” In Richard B. Braithwaite (ed.),  
9 *The Foundations of Mathematics and other Logical Essays*, 156–198, Routledge  
10 and Kegan Paul, London.
- 11 Reprinted in Henry E. Kyburg Jr. & Howard E. Smokler (1964, eds.) *Studies in  
12 Subjective Probability*, 61–92, Wiley, New York. (2<sup>nd</sup> edn. 1980, Krieger, New  
13 York.)
- 14 Rieger, Marc Oliver & Mei Wang (2008) “Prospect Theory for Continuous  
15 Distributions,” *Journal of Risk and Uncertainty* 36, 83–102.
- 16 Ruggeri, Kai, Sonia Alí, Mari Louise Berge, et al. (2020) “Replicating Patterns of  
17 Prospect Theory for Decision under Risk,” *Nature Human Behavior* 4, 622–633.
- 18 Samuelson, Paul A. (1959) “The St. Petersburg Paradox as a Divergent Double  
19 Limit,” *International Economic Review* 1, 31–37.
- 20 Schneider, Mark, Jonathan W. Leland, & Nathaniel T. Wilcox (2018) “Ambiguity  
21 Framed,” *Journal of Risk and Uncertainty* 57, 133–151.
- 22 Tversky, Amos & Daniel Kahneman (1981) “The Framing of Decisions and the  
23 Psychology of Choice,” *Science* 211, 453–458.
- 24 Tversky, Amos & Daniel Kahneman (1986) “Rational Choice and the Framing of  
25 Decisions,” *Journal of Business* 59, S251–S278.
- 26 Tversky, Amos & Daniel Kahneman (1992) “Advances in Prospect Theory:  
27 Cumulative Representation of Uncertainty,” *Journal of Risk and Uncertainty* 5,  
28 297–323.
- 29 Savage, Leonard J. (1971) “Elicitation of Personal Probabilities and Expectations,”  
30 *Journal of the American Statistical Association* 66, 783–801.
- 31 Smith, Vernon L. (1982) “Microeconomic Systems as an Experimental Science,”  
32 *American Economic Review* 72, 923–955.

- 1 Starmer, Chris (2000) “Developments in Non-Expected Utility Theory: The Hunt for  
2 a Descriptive Theory of Choice under Risk,” *Journal of Economic Literature* 38,  
3 332–382.
- 4 van de Kuilen, Gijs & Peter P. Wakker (2011) “The Midweight Method to Measure  
5 Attitudes toward Risk and Ambiguity,” *Management Science* 57, 582–598.
- 6 von Winterfeldt, Detlof & Ward Edwards (1982) “Costs and Payoffs in Perceptual  
7 Research,” *Psychological Bulletin* 91, 609–622.
- 8 Wakker, Peter P. (2010) “*Prospect Theory: For Risk and Ambiguity.*” Cambridge  
9 University Press, Cambridge, UK.
- 10 Wakker, Peter P. (2020), “Annotated Bibliography.”  
11 <http://personal.eur.nl/wakker/refs/webrfrncs.docx>
- 12 Wakker, Peter P., Ido Erev, & Elke U. Weber (1994) “Comonotonic Independence:  
13 The Critical Test between Classical and Rank-Dependent Utility Theories,”  
14 *Journal of Risk and Uncertainty* 9, 195–230.
- 15 Weber, Elke U. & Britt Kirsner (1997) “Reasons for Rank-Dependent Utility  
16 Evaluation,” *Journal of Risk and Uncertainty* 14, 41–61.
- 17 Wilcox, Nathaniel T. (1993) “Lottery Choice: Incentives, Complexity and Decision  
18 Time,” *Economic Journal* 103, 1397–1417.
- 19 Wu, George (1994) “An Empirical Test of Ordinal Independence,” *Journal of Risk  
20 and Uncertainty* 9, 39–60.
- 21

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

**ONLINE APPENDIX OF**

**A DEFENSE OF PROSPECT THEORY IN BERNHEIM &  
SPRENGER’S EXPERIMENT: A REVIEW OF  
COMPLEXITY AVERSION**

MOHAMMED ABDELLAOUI<sup>a</sup>, CHEN LI<sup>b</sup>, PETER P. WAKKER<sup>b</sup>, & GEORGE WU<sup>c</sup>

a: HEC Paris and CNRS, Jouy-en-Josas, France, Abdellaoui@hec.fr

b: Erasmus School of Economics, Erasmus University Rotterdam, Rotterdam, the Netherlands, c.li@ese.eur.nl; Wakker@ese.eur.nl

c: Booth School of Business, University of Chicago, Chicago, USA, wu@chicagobooth.edu

20 August, 2020

Bernheim & Sprenger (2020), BS henceforth, suggest that people are usually complexity averse. Complexity here refers only to the number of outcomes of lotteries. Aversion to more comprehensive or different forms of complexity of lotteries has been studied by Armantier & Treich (2016), Bruce & Johnson (1996), Kovarik, Levin, & Wang (2016), Mador, Sonsino, & Benzion (2000), and Sonsino, Benzion, & Mador (2002). Following BS, we focus only on empirical studies that have investigated the number of outcomes. BS cite five papers on complexity aversion in their footnote 70. The first four, Iyengar & Kamenica (2010), Iyengar & Lepper (2000), Iyengar, Jiang, & Huberman (2003), and Sonsino & Mandelbaum (2001), considered a different topic, preference against flexibility (number of available choice options to choose one from). The fifth, Stodder (1997), is on confusions of averages versus marginals and complexity of multiple stage lotteries, which, again, are different topics. (It is also only theoretical, with no data.) Hence, these references will not be considered here.

Complexity aversion here means that, other things equal, people prefer lotteries with few outcomes because they are less complex. We review the literature, focusing on gains, the domain considered by BS. We conclude that the prevailing finding is the opposite. That is, people usually prefer lotteries with many outcomes to few outcomes, and, in this sense, they are complexity seeking.

1       BS find that certainty equivalents of lotteries (0.4: 30, 0.6: 20) exceed those of  
2 (0.4: 30, 0.3: 20 +  $\varepsilon$ , 0.3: 20 -  $\varepsilon$ ) considerably, even for small  $\varepsilon > 0$ , with similar  
3 findings for (0.6: 30, 0.4: 20) versus (0.3: 30 +  $\varepsilon$ , 0.3: 30 -  $\varepsilon$ , 0.4: 20). By any  
4 rational theory, the certainty equivalents should, to the contrary, be almost the same  
5 for small  $\varepsilon$ . On the basis of these two observations, BS conclude that people are  
6 complexity averse.

7       Many studies have tested special preferences for numbers of outcomes, usually  
8 considering a pure case: certainty equivalents are measured for different framings of  
9 identical lotteries, for instance (0.4: 30, 0.6: 20) versus (0.4: 30, 0.3: 20 0.3: 20).  
10 Although all rational theories of choice require identical certainty equivalents,  
11 experiments find systematic violations. Here a pure effect of perceived number of  
12 outcomes occurs. Terms used to designate such violations include boundary effects,  
13 violations of coalescing (collapsing), and event/outcome splitting effects. The latter  
14 term is sometimes (e.g., in works by Humphrey, Starmer, and Sugden) combined with  
15 a directional assumption, being complexity seeking. Birnbaum does not add this  
16 directional assumption to this term.

17       Violations of coalescing can be taken as a special case of the attribute splitting  
18 effect (Weber, Eisenführ, & von Winterfeldt 1988), or the part-whole bias (Bateman  
19 et al. 1997), or the unpacking effect (Tversky & Koehler 1994). Here splitting  
20 something up increases the total weight. This underlies Birnbaum's theories. He  
21 studied violations of coalescing most extensively. His RAM and TAX models predict  
22 that splitting the best outcome of a lottery improves the lottery, but splitting the worst  
23 outcome (also if a gain) worsens it. If one normalizes decision weights to always add  
24 to 1, as in Birnbaum's models, then Birnbaum's predictions will hold. Then increasing  
25 the weight of the worst outcome indeed worsens the value. If one does not normalize  
26 the weights, as in separable prospect theory, then increasing the weights of gains  
27 (whether best or worst) improves the value. Combining these ideas suggests a strong  
28 preference for event splitting if it concerns the best outcome, and less clear effects for  
29 the worst outcome, but probably a preference against. Overall, we can then expect  
30 more preference for than against event splitting. In other words, the preceding  
31 arguments suggest more complexity seeking than aversion. This is indeed what our  
32 literature review finds.



1 Our literature search is based on searching the terms “boundary,” “collaps,”  
 2 “coalesc,” “complex,” and “split” in Wakker (2020), where we excluded the certainty  
 3 effect and followed up on cited papers.

- 4
- 5 • The following three papers report prevailing complexity aversion:  
 6 Bernheim & Sprenger (2020), Huck & Weizsäcker (1999), Moffatt, Sitzia, &  
 7 Zizzo (2015).

- 8
- 9 • The following seven papers report prevailing complexity seeking:  
 10 Birnbaum (2005), Birnbaum (2007), Humphrey (1995), Humphrey (2000),  
 11 Humphrey (2001a), Humphrey (2006), Starmer & Sugden (1993).

- 12
- 13 • The following five papers report about as much aversion as seeking:  
 14 Birnbaum (2004), Birnbaum, Schmidt, & Schneider (2017), Schmidt & Seidl  
 15 (2014), Humphrey (2001b), Weber (2007).

16

17 We conclude that the findings on complexity aversion are volatile, but the literature  
 18 has documented more complexity seeking than aversion for gains.

19

## 20 REFERENCES FOR ONLINE APPENDIX

- 21 Armantier, Olivier & Nicolas Treich (2016) “The Rich Domain of Risk,”  
 22 *Management Science* 62, 1954–1969.
- 23 Bateman, Ian J., Alistair Munro, Bruce Rhodes, Chris Starmer, & Robert Sugden  
 24 (1997) “Does Part-Whole Bias Exist? An Experimental Investigation,” *Economic*  
 25 *Journal* 107, 322–332.
- 26 Bernheim, B. Douglas & Charles Sprenger (2020) “On the Empirical Validity of  
 27 Cumulative Prospect Theory: Experimental Evidence of Rank-Independent  
 28 Probability Weighting,” *Econometrica* 88, 1363–1409.
- 29 Birnbaum, Michael H. (2004) “Causes of Allais Common Consequence Paradoxes:  
 30 An Experimental Dissection,” *Journal of Mathematical Psychology* 48, 87–106.
- 31 Birnbaum, Michael H. (2005) “Three New Tests of Independence That Differentiate  
 32 Models of Risky Decision Making,” *Management Science* 51, 1346–1358.

- 1 Birnbaum, Michael H. (2007) “Tests of Branch Splitting and Branch-Splitting  
2 Independence in Allais Paradoxes with Positive and Mixed Consequences,”  
3 *Organizational Behavior and Human Decision Processes* 102, 154–173.
- 4 Birnbaum, Michael H., Ulrich Schmidt, & Miriam D. Schneider (2017) “Testing  
5 Independence Conditions in the Presence of Errors and Splitting Effects,” *Journal*  
6 *of Risk and Uncertainty* 54, 61–85.
- 7 Bruce, Alistair C. & Johnnie E. V. Johnson (1996) “Decision-Making under Risk: The  
8 Effect of Complexity on Performance,” *Psychological Reports* 79, 67–76.
- 9 Huck, Steffen & Georg Weizsäcker (1999) “Risk, Complexity, and Deviations from  
10 Expected-Value Maximization: Results of a Lottery Choice Experiment,” *Journal*  
11 *of Economic Psychology* 20, 699–715.
- 12 Humphrey, Steven J. (1995) “Regret Aversion or Event-Splitting Effects? More  
13 Evidence under Risk and Uncertainty,” *Journal of Risk and Uncertainty* 11, 263–  
14 274.
- 15 Humphrey, Steven J. (2000) “The Common Consequence Effect: Testing a Unified  
16 Explanation of Recent Mixed Evidence,” *Journal of Economic Behavior and*  
17 *Organization* 41, 239–263.
- 18 Humphrey, Steven J. (2001a) “Are Event-Splitting Effects Actually Boundary  
19 Effects?,” *Journal of Risk and Uncertainty* 22, 79–93.
- 20 Humphrey, Steven J. (2001b) “Non-transitive Choice: Event-Splitting Effects or  
21 Framing Effects?,” *Economica* 68, 77–96.
- 22 Humphrey, Steven J. (2006) “Does Learning Diminish Violations of Independence,  
23 Coalescing and Monotonicity?,” *Theory and Decision* 61, 93–128.
- 24 Iyengar, Sheena S., & Emir Kamenica (2010) “Choice Proliferation, Simplicity  
25 Seeking, and Asset Allocation,” *Journal of Public Economics* 94, 530–539.
- 26 Iyengar, Sheena S., & Mark R. Lepper (2000) “When Choice Is Demotivating: Can  
27 One Desire Too Much of a Good Thing?” *Journal of Personality and Social*  
28 *Psychology* 79, 995–1006.
- 29 Iyengar, Sheena S., Wei Jiang, & Gur Huberman (2003) “How Much Choice Is Too  
30 Much: Contributions to 401(k) Retirement Plans,” Pension Research Council  
31 Working Paper 2003-10.
- 32 Kovarik, Jaromir, Dan Levin & Tao Wang (2016) “Ellsberg Paradox: Ambiguity and  
33 Complexity Aversions Compared,” *Journal of Risk and Uncertainty* 52, 47–64.

- 1 Mador, Galit, Doron Sonsino, & Uri Benzion (2000) “On Complexity and Lotteries’  
2 Evaluations — Three Experimental Observations,” *Journal of Economic*  
3 *Psychology* 21, 625–637.
- 4 Moffatt, Peter G., Stefania Sitzia, & Daniel John Zizzo (2015) “Heterogeneity in  
5 Preferences towards Complexity,” *Journal of Risk and Uncertainty* 51, 147–170.
- 6 Schmidt, Ulrich & Christian Seidl (2014) “Reconsidering the Common Ratio Effect:  
7 The Roles of Compound Independence, Reduction, and Coalescing,” *Theory and*  
8 *Decision* 77, 323–339.
- 9 Sonsino, Doron, Uri Benzion, & Galit Mador (2002) “The Complexity Effects on  
10 Choice with Uncertainty—Experimental Evidence,” *Economic Journal* 112, 936–  
11 965.
- 12 Sonsino, Doron, & Marvin Mandelbaum (2001) “On Preference for Flexibility and  
13 Complexity Aversion: Experimental Evidence,” *Theory and Decision* 51, 197–  
14 216.
- 15 Starmer, Chris & Robert Sugden (1993) “Testing for Juxtaposition and Event-  
16 Splitting Effects,” *Journal of Risk and Uncertainty* 6, 235–254.
- 17 Stodder, James (1997) “Complexity Aversion: Simplification in the Herrnstein and  
18 Allais Behaviors,” *Eastern Economic Journal* 23, 1–15.
- 19 Tversky, Amos & Derek J. Koehler (1994) “Support Theory: A Nonextensional  
20 Representation of Subjective Probability,” *Psychological Review* 101, 547–567.
- 21 Wakker, Peter P. (2020), “Annotated Bibliography.”  
22 <http://personal.eur.nl/wakker/refs/webrfrncs.docx>
- 23 Weber, Bethany J. (2007) “The Effects of Losses and Event Splitting on the Allais  
24 Paradox. Judgment in Decision Making,” *Judgment and Decision Making* 2, 115–  
25 125.
- 26 Weber, Martin, Franz Eisenführ, & Detlof von Winterfeldt (1988) “The Effects of  
27 Splitting Attributes on Weights in Multiattribute Utility Measurement,”  
28 *Management Science* 34, 431–445.
- 29  
30  
31