

Correcting Proper Scoring Rules for Risk Attitudes¹

Theo Offerman^a, Joep Sonnemans^a, Gijs van de Kuilen^a, & Peter P. Wakker^b

a: CREED, Dept. of Economics, University of Amsterdam, Roetersstraat 11,
Amsterdam, 1018 WB, The Netherlands

b: Econometric Institute, Erasmus University, P.O. Box 1738, Rotterdam, 3000 DR, the
Netherlands

May, 2006

10 ABSTRACT. Proper scoring rules, introduced in the 1950s, efficiently elicit subjective beliefs
11 about likelihoods, but do so only under the assumption of expected value maximization. The
12 latter assumption can be violated because of nonlinear utility (Bernoulli), nonexpected utility
13 (Allais), and ambiguity attitudes for unknown probabilities (Keynes, Knight, Ellsberg). These
14 violations of expected value have been incorporated in modern decision theories. This paper
15 shows how proper scoring rules can be generalized to those modern theories, and can become
16 valid under risk aversion and other deviations from expected value. An experiment
17 demonstrates the feasibility of our extension, yielding plausible empirical results. Violations of
18 additivity of subjective probabilities are reduced by our extension, although they do not
19 disappear entirely, which suggests genuine nonadditivity in subjective beliefs. The quality of
20 reported probabilities is better under repeated small incentives than under single large
21 incentives.

23 Key words: belief measurement, proper scoring rules, ambiguity, Knightian uncertainty,
24 subjective probability, nonexpected utility

¹ A preliminary version of this paper circulated with the title “Is the Quadratic Scoring Rule Really Incentive Compatible?”

25 1. Introduction

26 In many situations, no probabilities are known of uncertain events that are relevant to our
 27 decisions, and subjective assessments of the likelihoods of such events have to be made.
 28 Proper scoring rules provide an efficient tool for eliciting such subjective assessments from
 29 choices. They use cleverly constructed optimization problems where the observation of one
 30 single choice suffices to determine the exact quantitative degree of belief of an agent. This
 31 procedure is more efficient than the observation of binary choices, as is most common in
 32 decision theory, because binary choices only give inequalities and approximations of beliefs.

33 The measurement of subjective beliefs is important in many domains (Gilboa &
 34 Schmeidler 1999; Manski 2004), and proper scoring rules have been widely used accordingly,
 35 in accounting (Wright 1988), Bayesian statistics (Savage 1971), business (Staël von Holstein
 36 1972), education (Echternacht 1972), medicine (Spiegelhalter 1986), psychology (Liberman
 37 & Tversky 1993; McClelland & Bolger 1994), and other fields (Hanson 2002, Prelec 2004).
 38 Proper scoring rules are especially useful for giving experts incentives to exactly specify their
 39 degrees of belief. They are commonly used, for instance, to improve the calibration of weather
 40 forecasters (Palmer & Hagedorn 2006). They have recently become popular in experimental
 41 economics and game theory. Advocates of the frequentist interpretation of probability can
 42 become more interested in subjective probabilities when exposed to proper scoring rules. The
 43 quadratic scoring rule is the most popular proper scoring rule today (McKelvey & Page 1990;
 44 Nyarko & Schotter 2002), and is the topic of this paper.

45 Proper scoring rules were introduced independently by Brier (1950), Good (1952, p. 112),
 46 and de Finetti (1962), and are based on the assumption of expected value maximization, i.e.
 47 risk neutrality. All applications up to today that we are aware of have maintained this
 48 assumption. Empirically, however, deviations from expected value maximization are common.
 49 First, Bernoulli (1738) pointed out that risk aversion prevails over expected value, so that,
 50 under expected utility, utility has to be concave rather than linear. Second, Allais (1953)
 51 demonstrated, for events with known probabilities, that people can be risk averse in ways that
 52 expected utility cannot accommodate, so that more general decision theories are called for with
 53 other factors in addition to utility curvature (Kahneman & Tversky 1979, Machina 1982,
 54 Quiggin 1982, Tversky & Kahneman 1992). Third, Keynes (1921), Knight (1921), and
 55 Ellsberg (1961) demonstrated the importance of ambiguity for events with unknown

56 probabilities (“Knightian uncertainty”), where phenomena occur that are fundamentally
 57 different than those for known probabilities, adding to the descriptive failure of expected value
 58 and showing that there are further factors generating deviations from expected value. Gilboa
 59 (1987), Gilboa & Schmeidler (1989), Schmeidler (1989), and Tversky & Kahneman (1992)
 60 developed decision theories that incorporate ambiguity.

61 This paper updates proper scoring rules from the expected-value model as assumed in the
 62 1950s, when proper scoring rules were introduced, to the current state of the art in decision
 63 theory. Thus, on the one hand, we extend the applicability and validity of proper scoring rules
 64 to the empirical phenomena put forward by modern decision theorists. On the other hand, we
 65 extend the applicability of those modern decision theories by introducing the efficient
 66 measurement technique of proper scoring rules into those theories. We give quantitative
 67 assessments of the distortions of classical proper scoring rules caused by the empirical
 68 deviations from expected value, and show how such distortions can be corrected for by means
 69 of our generalized method.

70 Our correction technique can be interpreted as a new method of calibration (Keren 1991,
 71 Yates 1990) that does not need many repeated observations, unlike traditional calibration
 72 techniques (Clemen & Lichtendahl 2002, Leher 2001). An efficient aspect of our method is
 73 that we need not elicit entire risk attitudes of agents so as to correct for them. For instance,
 74 for nonexpected utility we need not go through an entire measurement of the utility and
 75 probability weighting functions to apply our correction. Instead, we can immediately infer
 76 the correction from a limited set of readily observable data. We demonstrate the feasibility of
 77 our method in an experiment where we measure the subjective beliefs of participants about the
 78 future performance of stocks after provision of information about past performance. The
 79 empirical findings will confirm the usefulness of our method.

80 The first part of this paper (Sections 2-5) is theoretical and the second part (Sections 6
 81 and further) is experimental. Section 2 provides the definition of proper scoring rules, and
 82 Section 3 provides the classical analysis of proper scoring rules, which assumes expected
 83 value maximization. A decision-theoretic analysis of the three main deviations from
 84 expected value maximization, and their effects on proper scoring rules, is in Section 4.
 85 Section 5 presents risk-corrections, serving to correct for the biases described in Section 4.
 86 Section 6 contains a simple example to illustrate how the theory developed in the first five
 87 sections can be applied empirically. It shows in particular that many decision-theoretic
 88 details, presented in the first part to justify our correction procedures, need not be studied
 89 when applying our method empirically. Readers interested in applying our method

90 empirically can skip most of Sections 3-5, reading only Section 3 up to Theorem 3.1 and
 91 Corollary 5.1. Methodological details of our experiment are described in Section 7. Section
 92 8 presents results regarding the biases that we correct for, and Section 9 some implications of
 93 the corrections of such biases. Discussions and conclusions are in Sections 10 and 11.
 94 Appendix A presents proofs and technical results, Appendix B surveys the implications of
 95 modern decision theories for our measurements, and Appendix C presents details of the
 96 experimental instructions.

97

98 **2. Proper Scoring Rules; Definitions**

99 Let E denote an event of which an agent is uncertain about whether or not it obtains,
 100 such as snow in Amsterdam in March 1932, whether a stock's value will decrease during the
 101 next half year, whether a ball randomly drawn from 20 numbered balls will have a number
 102 below 5, whether the 100th digit of π is 3, etc. The degree of uncertainty of an agent about E
 103 will obviously depend on the information that the agent has about E . Some agents may even
 104 know with certainty about some of the events. Most events will, however, be uncertain. For
 105 most uncertain events, no objective probabilities of occurrence are known, and our decisions
 106 have to be based on subjective assessments, consciously or not, of their likelihood.

107 Prospects designate event-contingent payments. We use the general notation $(E:x, y)$ for
 108 a *prospect* that yields outcome x if event E obtains and outcome y if E^c obtains, with E^c the
 109 *complementary event* not- E . The unit of payment for outcomes is one dollar. *Risk* concerns
 110 the case of known probabilities. Here, for a prospect $(E:x, y)$, the probability p of event E is
 111 known, and we can identify this prospect with a probability distribution $(p:x, y)$ over money,
 112 yielding x with probability p and y with probability $1-p$.

113 Several methods have been used in the literature to measure the subjective degree of
 114 belief of an agent in event E .

- 115 • In decision theory, binary choices have been used for this purpose. If the agent prefers
 116 the prospect $(E:1, 0)$ to receiving 0.10 for sure and also to receiving 0.20 for sure, but
 117 rather receives 0.30 for sure than the prospect, then under some assumptions (expected
 118 value; see later) we may conclude that the subjective probability of E is between 0.20 and
 119 0.30. In this way, many choices must be observed to infer levels of belief, and only
 120 approximations are obtained.

- We may seek to observe an indifference between $(E:1, 0)$ and a sure gain r , in which case we can conclude (again assuming expected value) that the degree of belief in E is exactly r . It is, however, hard to obtain measurements of indifferences. A Becker-DeGroot-Marschak mechanism can be used. Some drawbacks of this method are that it is fairly complex (Albers, Pope, Selten, & Vogt 2000, p. 116; Braga & Starmer 2005), needs further assumptions (Becker, de Groot, & Marschak 1964, p. 226; Guala 2000; Karni & Safra 1987), and is prone to irrational auction strategies (Plott & Zeiler 2005 p. 537).
- de Finetti (1931, 1937) devised a clever book-making technique with agents having to announce subjective probabilities, others taking them up on the corresponding betting odds, and the agents vulnerable to money pumps as soon as they violate the Bayesian laws of probability. This technique, however, concerns hypothetical decision situations and cannot be implemented in practice.
- Introspective questions can be used, simply asking what the degree of belief r is. It is, however, not clear why an agent would answer these questions truthfully and, without incentives, it may even be unclear to the agent (e.g. when the agent is a frequentist) what this number r is supposed to signify.

Proper scoring rules provide an efficient and operational manner for measuring subjective beliefs that deliver what the above methods seek to do, while avoiding the problems mentioned. They do not consider binary choice as in most decision theories but instead a multiple-choice optimization problem. They ask for the value r above without need to observe indifferences and with proper incentives added, and have been widely implemented unlike the book-making devise.

Under the *quadratic scoring rule (QSR)*, the most commonly used proper scoring rule and the rule considered in this paper, a *qsr-prospect*

$$(E: 1-(1-r)^2, 1-r^2), \quad (2.1)$$

is offered to the agent, where $0 \leq r \leq 1$ is a number that the agent can choose freely. The number chosen is a function of E , sometimes denoted r_E , and is called the (*uncorrected reported probability*) of E . The reasons for this term will be explained later. More general prospects $(E: a-b(1-r)^2, a-br^2)$ for any $b>0$ and $a \in \mathbb{R}$ can be considered, but for simplicity we restrict our attention to $a = b = 1$. No negative payments can occur, so that the agent never loses money. It is obvious that if the agent is certain that E will obtain, then he will maximize $1-(1-r)^2$, irrespective of $1-r^2$, and will choose $r=1$. Similarly, $r=0$ is chosen if E

153 will certainly not obtain. The choice of $r = 0.5$ gives a riskless prospect, yielding 0.75 with
 154 certainty. Increasing r increases the payment under E but decreases it under E^c . Under the
 155 event that happens, the QSR pays 1 minus the squared distance between the reported
 156 probability of a clairvoyant (who assigns probability 1 to the event that happens) and the
 157 reported probability of the agent (r under E , $1-r$ under E^c).

158 We, finally, discuss a symmetry between E and E^c , leading to a restriction on the data
 159 that can be obtained through QSRs and that will be further discussed in Subsection 4.3. For
 160 any prospect, the agent can have the payments under E and E^c interchanged by interchanging
 161 r and $1-r$. We display this point for future reference.

162

163 OBSERVATION 2.1. The quadratic scoring rule for event E presents the same choice of
 164 prospects as the quadratic scoring rule for event E^c , with each prospect resulting from r as
 165 reported probability of E identical to the prospect resulting from $1-r$ as reported probability
 166 of E^c . \square

167

168 Because of Observation 2.1, we have

$$169 \quad r_{E^c} = 1 - r_E. \quad (2.2)$$

170 This symmetry will hold throughout this paper, and will be crucial for the nonexpected utility
 171 analysis presented later.

172

173 **3. Proper Scoring Rules and Subjective Expected Value**

174 When proper scoring rules were introduced in the 1950s, they were based on the
 175 assumption of subjective expected value maximization. This assumption is still made in
 176 applications of proper scoring rules today, and will also be made in this section. It means,
 177 first, that the agent assigns a subjective probability p to each event E .² Under common

² In this paper, the term *subjective probability* is used only for probability judgments that are Bayesian in the sense of satisfying the laws of probability. In the literature, the term subjective probability has sometimes been used for judgments that deviate from the laws of probability, including cases where these judgments are nonlinear transformations of objective probabilities when the latter are given. Such concepts, different than

178 richness assumptions, the subjective probabilities must agree with objective probabilities
 179 whenever the latter are given to the agent³, and such agreement is assumed throughout this
 180 paper. Second, the agent maximizes expected value with respect to probabilities, irrespective
 181 of whether these are objective or subjective. These assumptions together constitute the
 182 *subjective expected value* model.

183 For QSRs and an event E with probability p, subjective expected value implies that the
 184 agent maximizes

$$185 \quad p \times (1 - (1-r)^2) + (1-p) \times (1-r^2) = 1 - p(1-r)^2 - (1-p)r^2. \quad (3.1)$$

186 If event E has probability p, then we also write R(p) for r_E throughout this paper. According
 187 to Eq. 3.1, and all other models considered in this paper, all events E with the same
 188 probability p have the same value r_E , so that R(p) is well-defined. We have the following
 189 corollary of Eq. 2.2.

$$190 \quad R(1-p) = 1 - R(p). \quad (3.2)$$

191 The following theorem demonstrates that the QSR is incentive compatible. The theorem
 192 immediately follows from the first-order optimality condition $2p(1-r) - 2r(1-p) = 0$ in Eq.
 193 3.1. Second-order optimality conditions are verified throughout this paper and will not be
 194 mentioned in what follows.

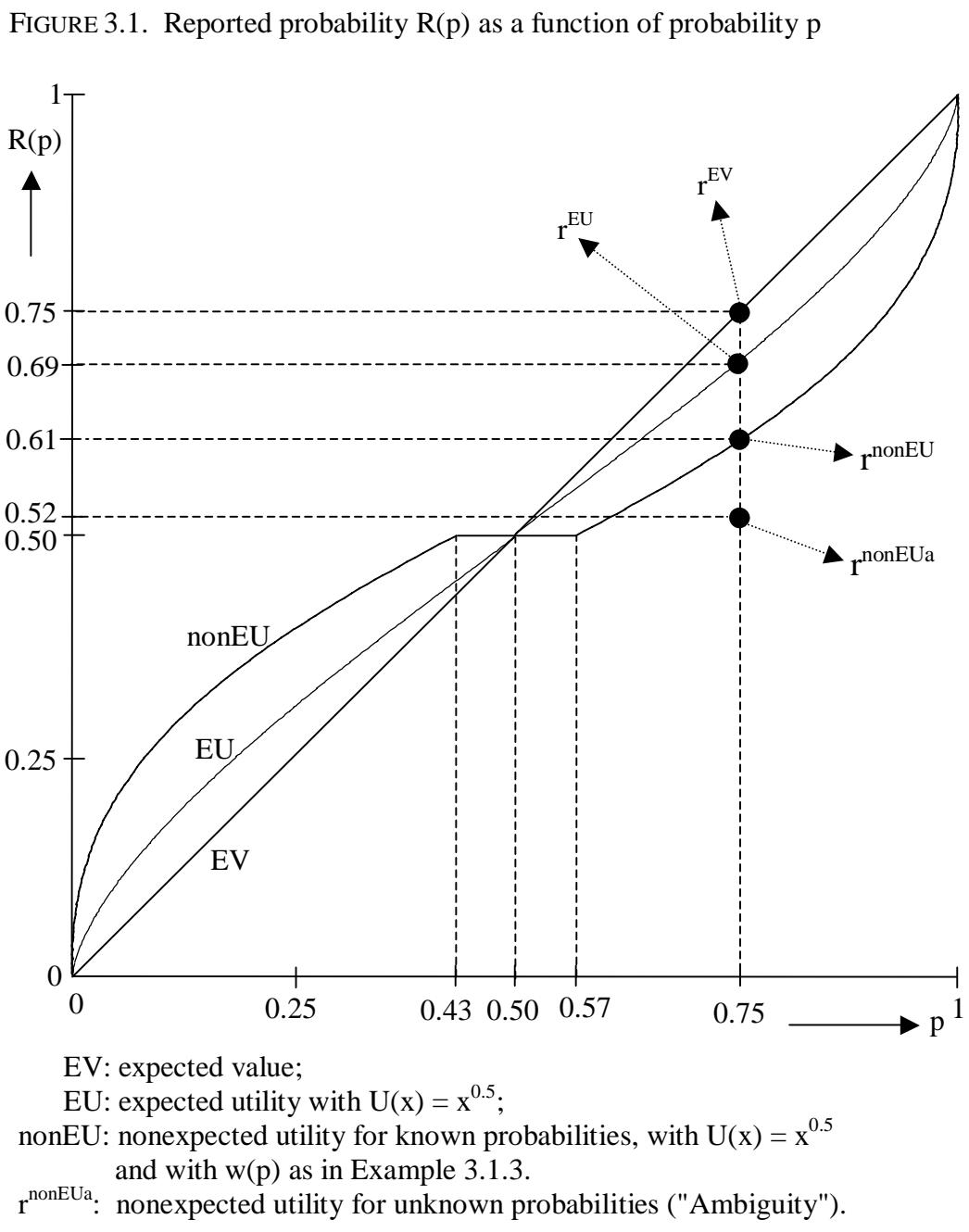
195
 196 THEOREM 3.1. Under subjective expected value maximization, the optimal choice r_E is equal
 197 to the (subjective or objective) probability p of event E, i.e. $R(p) = p$. \square
 198

199 It is in the agent's best interest to truly report his subjective probability of E. This
 200 explains the term "reported probability." Reported probabilities satisfy the Bayesian
 201 additivity condition for probabilities. We call the number r_E the (*uncorrected*) *reported*
 202 *probability*.

probabilities, will be analyzed in later sections, and we will use the term (probability) weights or beliefs,
 depending on the way of generalization, to designate them.

³ This follows first for equally-probable n-fold partitions of the universal event, where because of symmetry all events must have both objective and subjective probabilities equal to 1/n. Then it follows for all events with rational probabilities because they are unions of the former events. Finally, it follows for all remaining events by proper continuity or monotonicity conditions. There have been several misunderstandings about this point, especially in the psychological literature (Edwards 1954, p. 396; Schoemaker 1982, Table 1).

203 Figure 3.1 depicts $R(p)$ as a function of the probability p which, under expected value as
 204 considered here, is simply the diagonal $r = p$, indicated through the letters EV.
 205 The other
 206 curves and points in the figure will be explained later. Throughout the first part of this paper,
 we use variations of the following theoretical example.



232
 233 EXAMPLE 3.2. An urn K ("known" distribution) contains 25 Crimson, 25 Green, 25 Silver,
 234 and 25 Yellow balls. One ball will be drawn at random. C designates the event of a crimson
 235 ball drawn, and G, S, and Y are similar. E is the event that the color is not crimson, i.e. it is
 236 the event $C^c = \{G, S, Y\}$. Under expected value maximization, $r_E = R(0.75) = 0.75$ is optimal

237 in Eq. 2.1, yielding prospect $(E:0.9375, 0.4375)$ with expected value 0.8125. The point r_E is
 238 depicted as r^{EV} in Figure 3.1. Theorem 3.1 implies that $r_G = r_S = r_Y = 0.25$. We have $r_G + r_S$
 239 $+ r_Y = r_E$, and the reported probabilities satisfy additivity. \square

240

241 **4. Three Commonly Found Deviations from Subjective Expected
 242 Value and Their Implications for Quadratic Proper Scoring Rules**

243 This section describes three factors that generate deviations from expected value
 244 maximization and, hence, can distort the classical analyses of proper scoring rules. The
 245 effects of each factor in this section are illustrated in Figure 3.1, explained later, and their
 246 quantitative size will be illustrated through extensions of Example 3.2. Subsection 4.1
 247 considers the first factor generating deviations, being nonlinear utility under expected utility.
 248 This section extends an earlier study of this factor by Winkler & Murphy (1970). Bliss &
 249 Panigirtzoglou (2004) also corrected estimations of probability distributions for utility
 250 curvature, but did not do it in the context of proper scoring rules. Subsection 4.2 considers
 251 the second factor, i.e. violations of expected utility for known probabilities. Subsection 4.3
 252 considers the third factor, i.e. ambiguity because of unknown probabilities.

253

254 *4.1. The First Deviation: Utility Curvature*

255

256 Bernoulli (1738) put forward the first deviation from expected value. Because of the so-
 257 called St. Petersburg paradox, Bernoulli proposed that people maximize expected utility with
 258 respect to a *utility function* U , which we assume continuously differentiable with positive
 259 derivative everywhere, implying strict increasingness. We assume throughout that $U(0) = 0$.
 260 Eq. 3.1 is now generalized to

261 $pU(1-(1-r)^2) + (1-p)U(1-r^2) . \quad (4.1.1)$

262 The first-order optimality condition for r , and a rearrangement of terms (as in the proof of
 263 Theorem 4.3.2), implies the following result. For $r \neq 0.5$, the theorem also follows as a
 264 corollary of Theorem 4.3.2 and Eq. 3.2.

265

266 THEOREM 4.1.1. Under expected utility with p the (subjective or objective) probability of
 267 event E , the optimal choice $r = R(p)$ satisfies:

$$268 \quad r = \frac{p}{p + (1-p) \frac{U'(1-r^2)}{U'(1-(1-r)^2)}} . \quad (4.1.2)$$

269 \square

270

271 A utility correction is imposed, based on the marginal-utility ratio at the two prizes of the
 272 qsr-prospect. It implies that r deviates from the objective probability p if the marginal
 273 utilities at the two prizes are different. For concave utility and $r > 0.5$, so that E is judged
 274 more probable than E^c and receives the highest payment, the marginal-utility ratio will
 275 exceed 1 and r will be lower, and closer to 0.5, than p . For $r < 0.5$, r will also be closer to 0.5
 276 than p , because $1-r$, the reported probability of E^c , is so too. It follows that risk aversion
 277 moves people in the direction of the riskless prospect of $r = 0.5$ (Kadane & Winkler 1988, p.
 278 359), a phenomenon confirmed empirically by Huck & Weiszäcker (2002) and Winkler
 279 (1967).

280 Figure 3.1 depicts an example of the function r under expected utility, indicated by the
 281 letters EU, and is similar to Figure 3 of Winkler & Murphy (1970). The decision-based
 282 distortion in the direction of 0.5 is opposite to the overconfidence (probability judgments too
 283 far from 0.5) found mostly in direct judgments of probability without real incentives (Fischer
 284 1982; Fischhoff, Slovic & Lichtenstein 1977), and found among experts seeking to
 285 distinguish themselves (Keren 1991, p. 224 and 252; the “expert bias”, Clemen & Rolle
 286 2001).

287

288 EXAMPLE 4.1.2. Consider Example 3.2, but assume expected utility with $U(x) = x^{0.5}$.
 289 Substitution of Eq. 4.1.2 (or Theorem 4.3.2 below) shows that $r_E = R(0.75) = 0.69$ is optimal,
 290 depicted as r^{EU} in Figure 3.1, and yielding prospect $(E:0.91, 0.52)$ with expected value
 291 0.8094. The extra risk aversion generated by concave U has led to a decrease of r_E by 0.06
 292 relative to Example 3.2, distorting the probability elicited, and generating an expected-value
 293 loss of $0.8125 - 0.8094 = 0.0031$. This amount can be interpreted as an uncertainty
 294 premium, designating the profit margin for an insurance company. The uncertainty premium
 295 will be larger for deviations from expected value considered later. By Eq. 2.2, $r_C = 0.31$, and
 296 by symmetry $r_G = r_S = r_Y = 0.31$ too. The reported probabilities violate additivity, because r_G

297 $r_S + r_Y = 0.93 > 0.69 = r_E$. This violation in the data reveals that expected value does not
 298 hold. \square

299

300 The above example illustrates how observations of reported probabilities can be used to
 301 directly reveal violations of additivity empirically. We are not aware of such tests in the
 302 literature. Under such violations, reported probabilities may not truly reveal beliefs but may
 303 be distorted by other factors. We will report empirical tests of additivity in Sections 6 and 9.

304

305 OBSERVATION 4.1.3. Under expected utility with probability measure P , $r_E = 0.5$ implies
 306 $P(E) = 0.5$. Conversely, $P(E) = 0.5$ implies $r_E = 0.5$ if risk aversion holds. Under risk
 307 seeking, $r_E \neq 0.5$ is possible if $P(E) = 0.5$. \square

308

309 Theorem 4.1.1 clarifies the distortions generated by nonlinear utility, but it does not
 310 provide an explicit expression of $R(p)$, i.e. r as a function of p , or vice versa. It seems to be
 311 impossible, in general, to obtain an explicit expression of $R(p)$. We can, however, obtain an
 312 explicit expression of the inverse of $R(p)$, i.e. p in terms of r . For numerical purposes, $R(p)$
 313 can then be obtained as the inverse of that function—this is what we did in our numerical
 314 analyses, and how we drew Figure 3.1. The following result follows from algebraic
 315 manipulations or, for $r \neq 0.5$, as a corollary of Corollary 4.2.5 hereafter.

316

317 COROLLARY 4.1.4. Under Eq. 4.1.2, the optimal choice $r = R(p)$ satisfies:

$$318 \quad p = \frac{r}{r + (1-r) \frac{U'(1-(1-r)^2)}{U'(1-r^2)}} . \quad (4.1.3)$$

319 \square

320 The result shows that the relation between p and r is nonlinear, so that r will violate
 321 additivity, as soon as marginal utility U' is nonsymmetric about 0.5, which holds for all
 322 regular nonlinear utility functions. Hence, additivity of reported probability provides a
 323 critical test for linearity of utility under expected utility.

324

325 4.2. The Second Deviation: Nonexpected Utility for Known Probabilities
326

327 In the nonexpected utility analyses that follow, we will often restrict our attention to $r \geq$
 328 0.5. Results for $r < 0.5$ then follow by interchanging E and E^c , and the symmetry of
 329 Observation 2.1 and Eq. 2.2.

330 We say that event A is (*revealed*) *more likely than* event B if, for some positive outcome
 331 x , say $x = 100$, the agent prefers $(A:x, 0)$ to $(B:x, 0)$. In all models considered hereafter, this
 332 observation is independent of the outcome $x > 0$. In view of the symmetry of QSRs in
 333 Observation 2.1, for $r \neq 0.5$ the agent will always allocate the highest payment to the most
 334 likely of E and E^c . It leads to the following restriction of QSRs.

335

336 OBSERVATION 4.2.1. Under the QSR in Eq. 2.1, the highest outcome is always associated
 337 with the most likely event of E and E^c . \square

338

339 Hence, QSRs do not give observations about most likely events when endowed with the
 340 worst outcome. Similar restrictions apply to logarithmic proper scoring rules, as well as all
 341 other proper scoring rules as they have been applied in the literature so far.

342 Some details on weak inequalities and corner solutions are as follows. A choice of $r =$
 343 0.5 may be driven by risk aversion, so that no likelihood ordering between E and E^c can be
 344 concluded then. A choice of $r \neq 0.5$ (if close to 0.5), may be driven by risk seeking with
 345 equal likelihood of E and E^c . Only interior solutions with a strict inequality $r > 0.5$ combined
 346 with E being strictly less likely than E^c are excluded for QSRs.

347 We now turn to the second deviation from de Finetti's assumption of expected value, put
 348 forward by Allais (1953), which deviates from Bernoulli's expected utility, and still pertains
 349 to events E with known probability p . With M denoting 10^6 , the preferences $M >$
 350 $(0.8: 5M, 0)$ and $(0.25:M, 0) < (0.20:5M, 0)$ are plausible. They would imply, under
 351 expected utility with $U(0) = 0$, the contradictory inequalities $U(M) > 0.8 \times U(5M)$ and
 352 $0.25U(M) < 0.20 \times U(5M)$ (implying $U(M) < 0.8 \times U(5M)$), so that they falsify expected
 353 utility. It has since been shown that this paradox does not concern an exceptional
 354 phenomenon pertaining only to hypothetical laboratory choices with extreme amounts of
 355 money, but that the phenomenon is relevant to real decisions for realistic stakes (Kahneman
 356 & Tversky 1979). The Allais paradox and other violations of expected utility have led to
 357 several alternative models for decision under risk, the so-called nonexpected utility models.

358 For the prospects relevant to this paper, QSRs with only two outcomes and no losses, all
 359 currently popular nonexpected-utility evaluations of qsr-prospects (Eq. 2.1) are of the
 360 following form (see Appendix B). We first present such evaluations for the case of highest
 361 payment under event E, i.e. $r \geq 0.5$, which can be combined with $p \geq 0.5$.

362 For $r \geq 0.5$: $w(p)U(1-(1-r)^2) + (1-w(p))U(1-r^2)$. (4.2.1)

363 Here w is a continuous strictly increasing function with $w(0) = 0$ and $w(1) = 1$, and is called a
 364 *probability weighting function*. Expected utility is the special case of $w(p) = p$. By
 365 symmetry, the case $r < 0.5$ corresponds with a reported probability $1-r > 0.5$ for E^c , giving
 366 the following representation.

367 For $r < 0.5$: $w(1-p)U(1-r^2) + (1 - w(1-p))U(1-(1-r)^2)$. (4.2.2)

368 The different weighting of an event when it has the highest or lowest outcome is called rank-
 369 dependence. It suffices, by Eqs. 2.2 and 3.2, to analyze the case of $r \geq 0.5$ for all events.

370 Both in Eq. 4.2.1 and in Eq. 4.2.2, w is applied only to probabilities $p \geq 0.5$, and needs to
 371 be assessed only on this domain in what follows. This restriction is caused by Observation
 372 4.2.1. We display the implication.

373

374 OBSERVATION 4.2.2. For the QSR, only the restriction of w to $[0.5,1]$ plays a role, and w 's
 375 behavior on $[0,0.5]$ is irrelevant. \square

376

377 Hence, for the risk-correction introduced later, we need to estimate w only on $[0.5,1]$. An
 378 advantage of this point is that the empirical findings about w are uncontroversial on this
 379 domain, the general finding being that w underweights probabilities there.⁴ Under
 380 nonexpected utility, not only a utility correction must be imposed, but also a probability
 381 weighting correction w must be applied to p , leading to the following result.

382

383 THEOREM 4.2.3. Under nonexpected utility with p the probability of event E, the optimal
 384 choice $r = R(p)$ satisfies:

⁴ On $[0,0.5]$ the patterns is less clear, with both underweighting and overweighting (Abdellaoui 2000, Bleichrodt & Pinto 2000, Gonzalez & Wu 1999).

385 For $r > 0.5$: $r = \frac{w(p)}{w(p) + (1-w(p))\frac{U'(1-r^2)}{U'(1-(1-r)^2)}} . \quad (4.2.3)$

386 □

387

The above result, again, follows from the first-order optimality condition, and also follows as a corollary of Theorem 4.3.2 below. As an aside, the theorem shows that QSRs provide an efficient manner for measuring probability weighting on $(0.5, 1]$ if utility is linear, because then simply $r = R(p) = w(p)$. An extension to $[0, 0.5]$ can be obtained by a modification of QSRs, discussed further in the next subsection (Eqs. 4.3.3 and 4.3.4).

393

EXAMPLE 4.2.4. Consider Example 4.1.2, but assume nonexpected utility with $U(x) = x^{0.5}$ and

$$396 \quad w(p) = \left(\exp(-(-\ln(p))^\alpha) \right) \quad (4.2.4)$$

397 (Prelec 1998), with parameter $\alpha = 0.65$. This function agrees with common empirical
 398 findings (Tversky & Kahneman 1992, Abdellaoui 2000, Bleichrodt & Pinto 2000, Gonzalez
 399 & Wu 1999). From Theorem 4.2.3 it follows that $r_E = R(0.75) = 0.61$ is now optimal,
 400 depicted as r^{nonEU} in Figure 3.1. It yields prospect $(E:0.85, 0.63)$ with expected value 0.7920.
 401 The extra risk aversion relative to Example 4.1.2 generated by w for this event E has led to an
 402 extra distortion of r_E by 0.08. The extra expected-value loss (uncertainty premium) relative
 403 to Example 4.1.2 is $0.8094 - 0.7920 = 0.0174$. By Eq. 4.2.1, $r_C = 0.39$, and by symmetry $r_G =$
 404 $r_S = r_Y = 0.39$ too. The reported probabilities strongly violate additivity, because $r_G + r_S + r_Y$
 405 $= 1.17 > 0.61 = r_E$. \square

406
407 The effects of probability weighting are strongest near $p = 0.75$ and, indeed, relative to
408 Example 4.1.2, nonexpected utility has generated a large extra distortion in the above
409 example. Figure 3.1 illustrates the effects through the curve indicated by nonEU. Note that
410 the curve is flat around $p = 0.5$, more precisely, on the probability interval $[0.43, 0.57]$. For
411 probabilities from this interval the risk aversion generated by nonexpected utility is so strong
412 that the agent goes for maximal safety and chooses $r = 0.5$, corresponding with the sure
413 outcome 0.75 (cf. Manski 2004 footnote 10). Such a degree of risk aversion is not possible
414 under expected utility, where $r = 0.5$ can happen only for $p = 0.5$ (Observation 4.1.3). This

415 observation cautions against assigning specific levels of belief to observations $r = 0.5$,
 416 because proper scoring rules may be insensitive to small changes in the neighborhood of $p =$
 417 0.5 . An explicit expression of p in terms of r , i.e. of $R^{-1}(p)$, follows next for $r > 0.5$,
 418 assuming that we can invert w .

419

420 COROLLARY 4.2.5. Under Eq. 4.2.1, the optimal choice $r = R(p)$ satisfies:

421 If $r > 0.5$, then $p = R^{-1}(r) = w^{-1}\left(\frac{r}{r + (1-r)\frac{U'(1-(1-r)^2)}{U'(1-r^2)}}\right)$. (4.2.5)

422 \square

423 In general, it may not be possible to derive both w and U from $R(p)$ without further
 424 assumptions, i.e. U and w may be nonidentifiable. Under regular assumptions about U and
 425 w , however, they have some different implications. The main difference is that, if we assume
 426 that U is differentiable (as done throughout this paper) and concave, then a flat part of $R(p)$
 427 around 0.5 must be caused by w (Observation 4.1.3).

428 Up to this point, we considered deviations from expected value and Bayesianism at the
 429 level of decision attitude, and beliefs themselves were not yet affected. This will be different
 430 in the next subsection.

431

432 4.3. The Third Deviation: Nonadditive Beliefs and Ambiguity for Unknown Probabilities

433

434 This section considers events for which no probabilities are known. It is commonly
 435 assumed in applications of proper scoring rules that the agent then chooses (Bayesian)
 436 subjective probabilities $p = P(E)$ for such events E , satisfying the laws of probability, and
 437 evaluates prospects the same way for subjective probabilities as if these probabilities were
 438 objective. Such an approach to unknown probabilities, staying as close as possible to risk, is
 439 called *probabilistic sophistication* (Machina & Schmeidler 1992). In traditional applications
 440 of proper scoring rules it is further assumed that the agent satisfies expected utility for known
 441 probabilities, but probabilistic sophistication is more general and allows deviations like those
 442 in the preceding subsection for known probabilities. Probabilistic sophistication can be
 443 interpreted as a last attempt to at least maintain Bayesianism at the level of beliefs. As we
 444 will see next, however, it fails descriptively. Empirical findings, initiated by Ellsberg (1961),

445 have demonstrated that probabilistic sophistication is commonly violated empirically. The
 446 following example gives details. For another kind of violation see Marinacci (2002).

447

448 EXAMPLE 4.3.1 [violation of probabilistic sophistication]. Consider Example 4.2.4, but now
 449 there is an additional urn A (“ambiguous”). Like urn K, A contains 100 balls colored
 450 Crimson, Green, Silver, or Yellow, but now the proportions of balls with these colors are
 451 unknown. C_a designates the event of a crimson ball drawn from A, and G_a, S_a , and Y_a are
 452 similar. E_a is the event $C_a^c = \{G_a, S_a, Y_a\}$. If probabilities are assigned to drawings from the
 453 urn A (as assumed by probabilistic sophistication) then, in view of symmetry, we must have
 454 $P(C_a) = P(G_a) = P(S_a) = P(Y_a)$, so that these probabilities must be 0.25. Then $P(E_a)$ must be
 455 0.75, as was $P(E)$ in Example 4.2.4. Under probabilistic sophistication combined with
 456 nonexpected utility as in Example 4.2.4, r_{E_a} must be the same as r_E in Example 4.2.4 for the
 457 known urn, i.e. $r_{E_a} = 0.61$. It implies that people must be indifferent between $(E:x, y)$ and
 458 $(E_a:x, y)$ for all x and y . The latter condition is typically violated empirically. People usually
 459 have a strict preference for known probabilities, i.e. $(E:x, y) > (E_a:x, y)$.⁵ Consequently, it is
 460 impossible to model beliefs about uncertain events E_a through probabilities, and probabilistic
 461 sophistication must fail. This observation also suggests that r_{E_a} may differ from r_E . \square

462

463 The deviations from expected value illustrated by the above example cannot be
 464 explained by utility curvature or probability weighting, and must be generated by other
 465 factors. Those other, new, factors refer to properties of beliefs and decision attitudes that are
 466 typical of unknown probabilities, and force us to give up on the additive measure $P(E)$ in our
 467 model. Besides decisions, also beliefs may deviate from the Bayesian principles.

468 The important difference between known and unknown probabilities was first
 469 emphasized by Keynes (1921) and Knight (1921). Keynes discussed the example of urns
 470 with unknown compositions. Demonstrations as in the above example, proving that it is
 471 impossible to account for observed behavior in terms of probabilities, were first given by
 472 Ellsberg (1961).

473 Studies of direct judgments have supported the thesis that subjective beliefs may deviate
 474 from Bayesian probabilities (Dempster 1968; McClelland & Bolger 1994; Shafer 1976;
 475 Tversky & Koehler 1994). Instead of a probability p of E , we have to substitute a general

⁵ This holds also if people can choose the three colors to gamble on in the ambiguous urn, so that there is no reason to suspect unfavorable compositions.

476 subjective function $B(E)$ in Eq. 4.2.1, with Eq. 4.2.2 adapted similarly, and with
 477 nonadditivity of B adding to the deviations from Bayesianism. B may be interpreted as a
 478 belief index, as it was in Schmeidler (1989) who initiated the use of nonadditive measures for
 479 unknown probabilities. There is no consensus in decision theory today about whether B can
 480 also comprise other components of decision attitude beyond beliefs.

481 In Example 4.3.1 we may have $B(C_a) = B(G_a) = B(S_a) = B(Y_a) < 0.25$, $B(E_a) < 0.75$, and
 482 $B(E_a) \neq B(G_a) + B(S_a) + B(Y_a)$. Such phenomena lead to the following general evaluation of
 483 the qsr-prospects of Eq. 2.1, for general events E .

484 For $r \geq 0.5$: $w(B(E))U(1-(1-r)^2) + (1-w(B(E)))U(1-r^2)$. (4.3.1)

485 The evaluation for $r < 0.5$ can again be obtained from Eq. 2.2. We give it for
 486 completeness.

487 For $r < 0.5$: $(1-w(B(E^c)))U(1-(1-r)^2) + w(B(E^c))U(1-r^2)$. (4.3.2)

488 In general, B assigns value 0 to the vacuous event \emptyset , value 1 to the universal event, and
 489 B is increasing in the sense that $C \supset D$ implies $B(C) \geq B(D)$. These properties obviously also
 490 hold for the composition $w(B(\cdot))$. This composition is called the *weighting function*, and is
 491 denoted W . In the literature, the weighting function W is usually taken as point of
 492 departure.⁶ Given W and strict increasingness of w , we can define $B = w^{-1}(W)$ so as to
 493 obtain consistency of notation.

494 The equality $B(E) + B(E^c) = 1$ (*binary additivity*) may very well be violated. Then it can
 495 be debated whether $B(E)$ or $1 - B(E^c)$, or some other index, is to be taken as index of belief,
 496 and whether other decision components beyond beliefs are comprised in some or all of these
 497 indexes. Such interpretations have not yet been settled, and further studies are called for.
 498 Whereas the interpretation of B as index of belief is open to debate, depending on further
 499 developments in decision theory, it is not open to debate that the behavioral component of
 500 risk attitude should be filtered out before an interpretation of belief can be considered.
 501 Filtering out this behavioral component, as a necessary preparation for further investigations
 502 of beliefs, is the contribution of this paper.

⁶Schmeidler (1989) and most other current studies of uncertainty assume, for simplicity, that w is the identity. Then B and W coincide.

503 As with the weighting function w under risk, B is also applied only to the most likely one
 504 of E and E^c in the above equations, reflecting again the restriction of the QSR of Observation
 505 4.2.1. Hence, under traditional QSR measurements we cannot test binary additivity directly
 506 because we measure $B(E)$ only when E is more likely than E^c . These problems can easily be
 507 amended by modifications of the QSR. For instance, we can consider prospects

508 $(E: 2-(1-r)^2, 1-r^2), \quad (4.3.3)$

509 i.e. qsr-prospects as in Eq. 2.1 but with a unit payment added under event E . The classical
 510 proper-scoring-rule properties of Section 2 are not affected by this modification, and the
 511 results of Section 3 are easily adapted. With this modification, we have the liberty to
 512 combine event E with the highest outcome both if E is more likely than E^c and if E is less
 513 likely, and we avoid the restriction of Observation 4.2.1. We then can observe w of the
 514 preceding subsection, and $W(E)$ and $B(E)$ over their entire domain. Similarly, with prospects

515 $(E: 1-(1-r)^2, 2-r^2), \quad (4.3.4)$

516 we can measure the duals $1 - W(E^c)$, $1 - w(1-p)$, and $1 - B(E^c)$ over their entire domain. In
 517 this study we confine our attention to the QSRs of Eq. 2.1 as they are classically applied
 518 throughout the literature, so as to reveal their biases according to the current state of the art of
 519 decision theory, suggesting remedies whenever possible, and signaling the problems that
 520 remain. We leave further investigations of the, we think promising, modifications of QSRs in
 521 the above equations to future studies.

522 The restrictions of the classical QSRs will also hold for the experiment reported later in
 523 this paper. There an application of the QSR to a small interval I is to be interpreted formally
 524 as the measurement of $1 - B(I^c)$. The restrictions also explain why the theorems below
 525 concern only the case of $r > 0.5$ (with $r = 0.5$ as a boundary solution).

526 The following theorem, our main theorem, specifies the first-order optimality condition
 527 for interior solutions of r for general decision making, incorporating all deviations described
 528 above.

529

530 THEOREM 4.3.2. Under Eq. 4.3.1, the optimal choice r satisfies:

531 If $r > 0.5$, then $r = r_E = \frac{w(B(E))}{w(B(E)) + (1-w(B(E)))\frac{U'(1-r^2)}{U'(1-(1-r)^2)}}.$ $(4.3.5)$

532 □

533

534 We cannot draw graphs as in Figure 3.1 for unknown probabilities, because the x-axis
 535 now concerns events and not numbers. The W values of ambiguous events will be relatively
 536 low for an agent with a general aversion to ambiguity, so that the reported probabilities r in
 537 Eq. 4.3.5 will be relatively small, i.e. close to 0.5. We give a numerical example.

538

539 EXAMPLE 4.3.3. Consider Example 4.3.1. Commonly found preferences $(E:100, 0) >$
 540 $(E_a:100, 0)$ imply that $w(B(E_a)) < w(B(E)) = w(0.75)$. Hence, by Theorem 4.3.2, r_{E_a} will be
 541 smaller than r_E . Given the strong aversion to unknown probabilities that is often found
 542 empirically (Becker & Brownson 1964; Camerer & Weber 1992), we will assume that $r_{E_a} =$
 543 0.52. It is depicted as r^{nonEU_a} in Figure 3.1, and yields prospect $(E_a:0.77, 0.73)$ with expected
 544 value 0.7596. The extra preference for certainty relative to Example 4.2.4 generated by
 545 unknown probabilities for this event E_a has led to an extra distortion of r_{E_a} by $0.61 - 0.52 =$
 546 0.09. The extra expected-value loss relative to Example 4.2.4 is $0.7920 - 0.7596 = 0.0324$.
 547 This amount can be interpreted as the ambiguity-premium component of the total uncertainty
 548 premium. By Eq. 4.2.1, $r_C = 0.48$, and by symmetry $r_G = r_S = r_Y = 0.48$ too. The reported
 549 probabilities violate additivity to an extreme degree, because $r_G + r_S + r_Y = 1.44 > 0.52 = r_{E_a}$.
 550 The behavior of the agent is close to a categorical fifty-fifty evaluation, where all nontrivial
 551 uncertainties are weighted the same without discrimination.

552 The belief component $B(E_a)$ is estimated to be $w^{-1}(W(E_a)) = w^{-1}(0.52) = 0.62$. This
 553 value implies that B must violate additivity. Under additivity, we would have $B(C_a) = 1 -$
 554 $B(E_a) = 0.38$ and then, by symmetry, $B(G_a) = B(S_a) = B(Y_a) = 0.38$, so that $B(G_a) + B(S_a) +$
 555 $B(Y_a) = 3 \times 0.38 = 1.14$. This value should, however, equal $B\{G_a, S_a, Y_a\} = B(E_a)$ under
 556 additivity which is 0.62, leading to a contradiction. Hence, additivity must be violated.

557 Of the total deviation of $r_{E_a} = 0.52$ from 0.75, being 0.23, a part of $0.06 + 0.08 = 0.14$ is
 558 the result of deviations from risk neutrality that distorted the measurement of $B(E_a)$, and 0.09
 559 is the result of nonadditivity (ambiguity) of belief B . □

560

561 Theorem 4.3.2 is valid for virtually all models of decision under uncertainty and
 562 ambiguity currently known in the literature, because Eqs. 4.3.1 and 4.3.2 capture all these
 563 models (see Appendix B). Some qualitative observations are as follows. If U is linear, then r
 564 = $w(B(E))$ follows for all $w(B(E)) > 0.5$, providing a very tractable manner of measuring the

565 nonadditive decision-theory measure $W = w \circ B$. A corner solution $r = 0.5$ results for all
 566 $w(B(E)) \leq 0.5$ with also $w(B(E^c)) \leq 0.5$, so that the classical QSR has no discriminatory
 567 power for such events. For events with known probabilities, such corner solutions
 568 correspond with the flat part of the nonEU curve in Figure 3.1. For ambiguous events,
 569 ambiguity aversion will enhance the existence of such corner solutions.

570 Expected utility is the special case where $W(B(E)) = P(E)$ for a probability measure P , so
 571 that Eqs. 4.3.1 and 4.3.2 are the same and each applies to all r . Theorem 4.1.1 demonstrated
 572 that Eq. 4.3.5 then also holds for $r \leq 0.5$.

573 In applications of proper scoring rules, we usually first observe $r = r_E$ and then want to
 574 derive $B(E)$ from r . The following corollary gives an explicit expression. It illustrates once
 575 more how deviations from expected utility (w) and nonlinear utility (the marginal-utility
 576 ratio) distort the classical proper-scoring-rule assumption of $B(E) = r$.

577

578 COROLLARY 4.3.4. Under Eq. 4.3.1, the optimal choice $r = r_E$ satisfies:

579 If $r > 0.5$, then $B(E) = w^{-1}\left(\frac{r}{r + (1-r)\frac{U'(1-(1-r)^2)}{U'(1-r^2)}}\right)$. (4.3.6)

580 \square

581

582 5. Measuring Beliefs through Risk Corrections

583 One way to measure $B(E)$ is by eliciting $W(E)$ and the function w from choices under
 584 uncertainty and risk, after which we can set

$$585 B(E) = w^{-1}(W(E)). \quad (5.1)$$

586 In general, such revelations of w and W are laborious. The observed choices depend not only
 587 on w and W but also on the utility function U , so that complex multi-parameter estimations
 588 must be carried out (Tversky & Kahneman 1992, p. 311).

589 A second way to elicit $B(E)$ is by measuring the *canonical probability* p of event E ,
 590 defined through the equivalence

$$(p:x, y) \sim (E:x, y) \quad (5.2)$$

for some preset $x > y$, say $x = 100$ and $y = 0$. Then $w(B(E))(U(x) - U(y)) = w(p)(U(x) - U(y))$, and $B(E) = p$ follows. Wakker (2004) discussed the interpretation of Eqs. 5.1 and 5.2 as belief.

Canonical probabilities were commonly used in early decision analysis (Raiffa 1968, Section 5.3; Yates 1990 pp. 25-27) under the assumption of expected utility. A recent experimental measurement is in Holt (2005, Chapter 30), who also assumed expected utility. A practical difficulty is that the measurement of canonical probabilities requires the measurement of indifferences, and these are not easily inferred from choice. For example, Holt (2005) used the Becker-deGroot-Marschak mechanism, discussed above. Huck & Weiszäcker (2002) compared the QSR to the measurement of canonical probabilities and found that the former is more accurate.

A third way to correct reported probabilities is to collect, for each r_E , many events to which the agent assigned the same r_E value in the past and of which we now know whether or not they obtained. We then determine the relative frequency of these events, assuming that this can be taken as the true, objective, probability. We, thus, turn these events into unambiguous events. The distance between r_E and this relative frequency is an index of the (mis)calibration of the agent. This calibration technique has been studied in theoretical game theory (Sandroni, Smorodinsky, & Vohra 2003), and has been applied to weather forecasters (Murphy & Winkler 1974). It needs extensive data, which is especially difficult to obtain for rare events such as earthquakes. It needs further assumptions about the stability of distortions over time, and is hard to apply in experiments of limited durations. These drawbacks were pointed out by Clemen & Lichtendahl (2002), who proposed correction techniques for probability estimates in the spirit of our paper, but still based these on traditional calibration techniques. Our correction (“calibration”) technique is considerably more efficient than traditional ones.

We now introduce risk corrections that combine the advantages of measuring $B(E) = w^{-1}(W(E))$, of measuring canonical probabilities, and of calibrating reported probabilities relative to objective probabilities, while avoiding the problems described above, by benefiting from the efficiency of proper scoring rules. The QSR does entail a restriction of the observations regarding $B(E)$ to cases of E being more likely than E^c (Observation 4.2.1).

Note that the right-hand sides of Eqs. 4.2.5 and 4.3.6 are identical. Hence, if we find a p with the same r value as E, then we can, because of Eq. 4.2.5, immediately substitute p for

624 the right-hand side of Eq. 4.3.6, getting $B(E) = p$ without need to know the ingredients w and
 625 U of Eq. 4.3.6. This observation (to be combined with Eq. 2.2 for $r < 0.5$) implies the
 626 following corollary, which we display for its empirical importance.

627

628 COROLLARY 5.1. Under Eq. 4.3.1, for the optimal choice $r = r_E$, assume that $r > 0.5$. Then

629
$$B(E) = R^{-1}(r). \quad (5.3)$$

630 \square

631

632 This corollary is useful for empirical purposes. It is the only implication of our
 633 theoretical analysis that is needed for applications. We first infer the (for the participant)
 634 optimal $R(p)$ for a set of exogenously given probabilities p that is so dense (all values $p = j/20$
 635 for $j \geq 10$ in our experiment) that we obtain a sufficiently accurate estimation of R and R^{-1} .
 636 Then, for all uncertain events E more likely than their complement, we immediately derive
 637 $B(E)$ from the observed r_E through Eq. 5.3. Summarizing:

638 If for event E the participant reports probability $r_E = r$
 639 and for objective probability p the participant also reports probability $R(p) = r$
 640 then $B(E) = p$.

641 We, therefore, directly measure the curve $R(p)$ in Figure 3.1 empirically, and apply its
 642 inverse to r_E . For $r_E = 0.5$, $B(E)$ and the inverse p may not be uniquely determined because of
 643 the flat part of R_{nonEU} in Figure 3.1.

644 We call the function R^{-1} the *risk-correction* (for proper scoring rules), and $R^{-1}(r_E)$ the
 645 *risk-corrected probability*. This value is the canonical probability, obtained without having
 646 measured indifferences such as through the Becker-DeGroot-Marschak mechanism, without
 647 having measured U and w as in decision theory, and without having measured relative
 648 frequencies in many repeated observations of past events with the same reported probabilities
 649 as in calibrations. Obviously, if $R(p)$ does not deviate much from p , then no risk correction is
 650 needed. Then reported probabilities r directly reflect beliefs, and we have ensured that
 651 traditional analyses of QSRs give proper results.

652 The curves in Figure 3.1 can be reinterpreted as inverses of risk corrections. The
 653 examples illustrated there were based on risk averse decision attitudes, leading to
 654 conservative estimations moved in the direction of 0.5. Risk seeking will lead to the opposite
 655 effect, and will generate overly extreme reported probabilities, suggesting overconfidence.
 656 Obviously, if factors in the probability elicitation of the calibration part induce

657 overconfidence and risk seeking, then our risk correction will detect those biases and correct
 658 for them. If, after the risk correction, overconfidence is (still) present, then it cannot be due
 659 to risk seeking. It then is convincing that overconfidence is a genuine property of belief,
 660 irrespective of risk seeking.

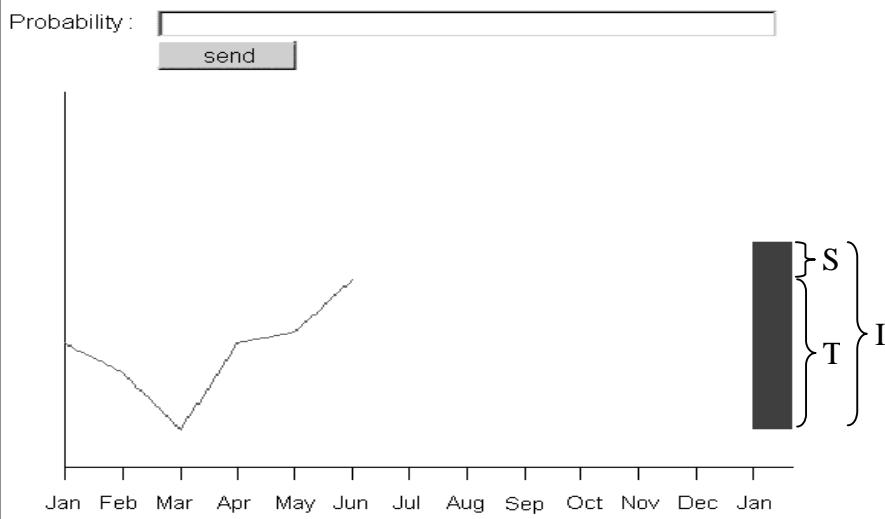
661 **6. An Illustration of Our Measurement of Belief**

662 This section describes risk corrections for a participant in the experiment so as to
 663 illustrate how our method can be applied empirically. It will illustrate that the conclusion in
 664 Corollary 5.1 is all of the theoretical analysis that is needed for applying our method. Results
 665 and curves for $r < 0.5$ are derived from $r > 0.5$ using Eq. 2.2; we will not mention this point
 666 explicitly in what follows.

667

668

669 **FIGURE 6.1. Layout of the screens**



680 The left side of Figure 6.1 displays the performance of stock 20 in our experiment from
 681 January 1 until June 1 1991 as given to the participants. It concerned CSM certificates
 682 dealing in sugar and bakery-ingredients. Further details (such as the absence of a unit on the
 683 y-axis) will be explained in Section 7. The right side of the figure displays two disjoint
 684 intervals S and T, and their union I = S ∪ T. For each of the intervals S, T, and I, participants
 685 reported the probability of the stock ending up in that interval on January 1 1992 (with some
 686 other questions in between these three questions). For participant 25, the results are as
 687 follows.

$$688 \quad r_s = 0.10; r_T = 0.55; r_I = 0.75. \quad (6.1)$$

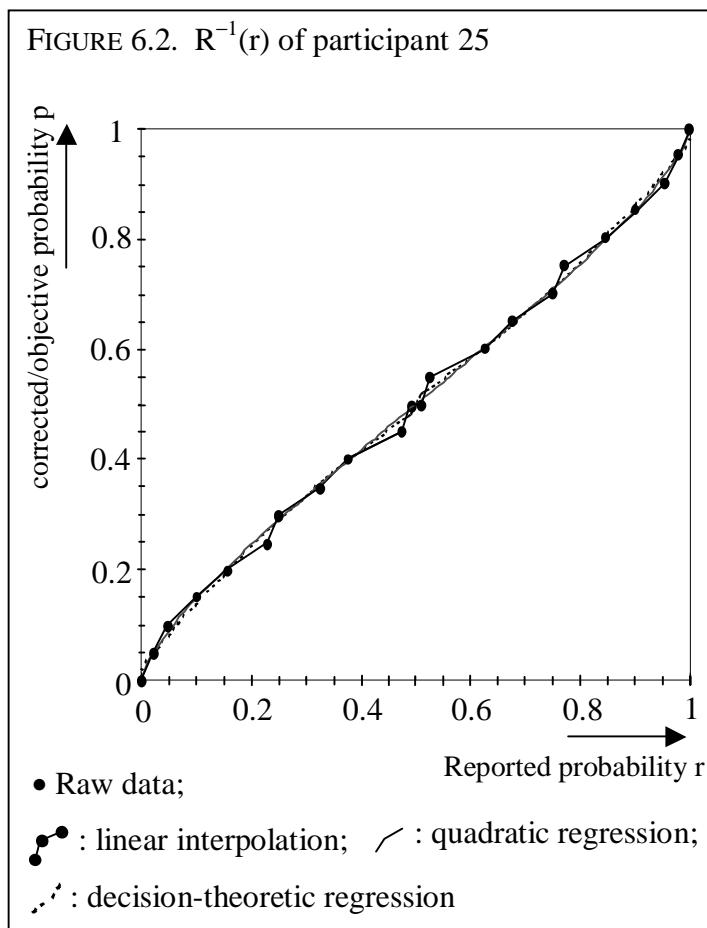
Under additivity of reported probability, $r_S + r_T - r_I$ (the *additivity bias*, defined in general in Eq. 7.5), should be 0, but here it is not and additivity is violated.

691 The additivity bias is $0.10 + 0.55 - 0.75 = -0.10$. (6.2)

Table 6.1 and Figure 6.2 (in inverted form) display the reported probabilities $R(p)$ that we measured from this participant, with the points in the figure indicating raw data, and the curves explained later.

TABLE 6.1. Reported probabilities $R(p)$ for given probabilities p of participant 25

p	.05	.10	.15	.20	.25	.30	.35	.40	.45	.50	.55	.60	.65	.70	.75	.80	.85	.90	.95
R(p)	.02	.04	.10	.15	.23	.25	.32	.37	.47	.51	.53	.63	.68	.75	.77	.85	.90	.96	.98



716 For simplicity of presentation, we analyze the data here using linear interpolation. Then
 717 $R(0.56)$ is obtained from

718 $0.56 = 0.8 \times 0.55 + 0.2 \times 0.60$, $R(0.55) = 0.53$, and $R(0.60) = 0.63$,

719 through

720 $R(0.56) = R(0.8 \times 0.55 + 0.2 \times 0.60) = 0.8 \times R(0.55) + 0.2 \times R(0.60) =$
 721 $0.8 \times 0.53 + 0.2 \times 0.63 = 0.55$.

722 Using this value for $R(0.56)$, further the values $R(0.15) = 0.10$ and $R(0.70) = 0.75$, and,
 723 finally, Eq. 5.3, we obtain the following risk-corrected beliefs.

724 $B(S) = R^{-1}(0.10) = 0.15$; $B(T) = R^{-1}(0.55) = 0.56$; $B(I) = R^{-1}(0.75) = 0.70$;
 725 the additivity bias is $0.15 + 0.56 - 0.70 = 0.01$. (6.3)

726 The risk-correction has reduced the violation of additivity, which according to Bayesian
 727 principles can be interpreted as a desirable move towards rationality. In the experiment
 728 described in the following sections we will see that this effect is statistically significant for
 729 single evaluations (treatment “t=ONE”), but is not so for repeated payments and decisions
 730 (treatment “t=ALL”).

731 It is statistically preferable to fit data with smoother curves than resulting from linear
 732 interpolation. We derived “decision-theoretic” parametric curves for $R(p)$ from Corollary
 733 4.2.5, with further assumptions explained at the end of Subsection 8.1.⁷ The resulting curve
 734 for participant 25 is given in the figure. $B = R^{-1}(r)$ and this curve lead to

735 $B(S) = R^{-1}(0.10) = 0.15$; $B(T) = R^{-1}(0.55) = 0.54$; $B(I) = R^{-1}(0.75) = 0.71$; the additivity
 736 bias is $0.15 + 0.54 - 0.71 = -0.02$, (6.4)

737 again reducing the uncorrected additivity bias. The quadratic curve gave virtually the same
 738 results, and will be discussed in Section 10.1.

739

⁷ The decision-theoretic curve in the figure is the function $p = B(E) = \frac{r}{r + (1-r)\frac{1.27(1-(1-r)^2)^{0.27}}{1.27(1-r^2)^{0.27}}}$, in

agreement with Corollaries 5.1 and 4.2.5, where we estimated $w(p) = p$ and found $\rho = 1.27$ as optimal value for $U(x)$ in Eq. 7.1.

7. An Experimental Application of Risk Corrections: Method

741 *Participants.* N = 93 students from a wide range of disciplines (45 economics; 13
 742 psychology, 35 other disciplines) participated in the experiment. They were self-selected
 743 from a mailing list of approximately 1100 people.

744

745 *Procedure.* Participants were seated in front of personal computers in 6 groups of
 746 approximately 16 participants each. They first received an explanation of the QSR, given in
 747 Appendix C. Then, for each uncertain event, participants could first report a probability (in
 748 percentages) by typing in an integer from 0 to 100. Subsequently, the confirmation screen
 749 displayed a list box with probabilities and the corresponding score when the event was (not)
 750 true, illustrated in Figure 7.1.

751

752

753 FIGURE 7.1.

754 Probability	755 Your score if statement is true	755 Your score if statement is not true
756 27%	4671	9271
28%	4816	9216
29%	4959	9159
30%	5100	9100
31%	5239	9039
32%	5376	8976
33%	5511	8911
34%	5644	8844
35%	5775	8775
36%	5904	8704

761 send

762

763
 764 All figures (including Figure 6.1) are reproduced here in black and white; in the experiment
 765 we used colors to further clarify the figures. The entered probability and the corresponding
 766 score were preselected in this list box. The participant could confirm the decision or change
 767 to another probability by using the up or down arrow or by scrolling to another probability
 768 using the mouse. The event itself was also visible on the confirmation screen. Thus, the
 769 reported probability r finally resulted for the uncertain event.

770

771 *Stimuli*

772 The participants provided 100 reported probabilities r for events with unknown probabilities
773 in the *stock-price part* of the experiment. For these events, we fixed June 1, 1991, as an
774 “evaluation date.” The uncertain events always concerned the question whether or not the
775 price of a stock would lie in a target-interval seven months after the evaluation date. For
776 each stock, the participants received a graph depicting the price of the stock on 0, 1, 2, 3, 4,
777 and 5 months prior to the evaluation date, as well as an upper and lower bound. Figure 6.1,
778 without the braces and letters, gives an example of the layout. We used 32 different stocks,
779 all real-world stock market data from the 1991 Amsterdam Stock Exchange. After 4 practice
780 questions, the graph of each stock-price was displayed once in the questions 5-36, once in the
781 questions 37-68, and once in the questions 69-100. We, thus, obtained three probabilistic
782 judgments of the performance of each stock, once for a large target-interval and twice for
783 small target-intervals that partitioned the large target-interval (see Figure 6.1). We partially
784 randomized the order of presentation of the elicitations. Each stock was presented at the
785 same place in the first, second, and third 32-tuple of elicitations, so as to ensure that
786 questions pertaining to the same stock were always far apart. The order of presentation of the
787 two small and one big interval within one stock were not randomized stochastically, but were
788 varied systematically, so that all orders of big and small intervals occurred equally often. We
789 also maximized the variation of whether small intervals were both very small, both
790 moderately small, or one very small and one moderately small.

791
792 In the *calibration part* of the experiment, participants made essentially the same decisions as
793 in the stock-price part, but now for 20 events with objective probabilities. Thus, participants
794 simply made choices between risky prospects with objective probabilities. We used two 10-
795 sided dice to determine the outcome of the different prospects and obtained measurements of
796 the reported probabilities corresponding to the objective probabilities 0.05, 0.10, 0.15, ...,
797 0.85, 0.90, and 0.95 (we measured the objective probability 0.95 twice). The event with
798 probability 0.25 was, for instance, described as “The outcome of the roll with two 10-sided
799 dice is in the range 01–25.” The decision screen was very similar to Figure 7.1, except for
800 the fact that we wrote “row-percentage” instead of “probability” and “your score if the roll of
801 the die is 01-25” instead of “your score if statement is true;” etc.

802
803 *Motivating participants.* Depending on whether the uncertain event obtained or not and on
804 the reported probability for the uncertain event, a number of points was determined for each

805 question through the QSR (Eq. 2.1), using 10000 points as unit of payment so as to have
 806 integer scores with four digits of precision. Thus, the maximum score for one question was
 807 10000, the minimum score was 0, and the certain score resulting from reported probability
 808 0.5 was 7500 points.

809 In treatment $t=ALL$, the sum of all points for all questions was calculated for each
 810 participant and converted to money through an exchange rate of 60000 points = €1, yielding
 811 an average payment of €15.05 per participant. For the calibration part we then used a box
 812 with twenty separate compartments containing pairs of 10-sided dice to determine the
 813 outcome of each of the twenty prospects at the same time for the treatment $t=ALL$.

814 In treatment $t=ONE$, the random-lottery incentive system was used. That is, at the end of
 815 the experiment, one out of the 120 questions that they answered was selected at random for
 816 each participant and the points obtained for this question were converted to money through
 817 an exchange rate of 500 points = €1, yielding an average payment of €15.30 per participant.

818 All payments were done privately at the end of the experiment.

819

820 *Analysis.* For the calibration part we only need to analyze probabilities of 0.5 or higher, by
 821 Observation 4.2.2. Indeed, by Eq. 3.2, every observation for $p < 0.5$ amounts to an
 822 observation for $p' = 1-p > 0.5$. It implies that we have two observations for all $p > 0.5$ (and
 823 three for $p = 0.95$).

824 We first analyze the data at the group level, assuming homogeneous participants. We
 825 start from the general model of Eq. 4.2.1. We used parametric fittings, where for U we used
 826 the *power utility with parameter ρ* , also known as the family of constant relative risk aversion
 827 (CRRA)⁸, and the most popular parametric family for fitting utility, which is defined as
 828 follows:

$$\begin{aligned} 829 \quad & \text{For } \rho > 0: U(x) = x^\rho; \\ 830 \quad & \text{for } \rho = 0: U(x) = \ln(x); \\ 831 \quad & \text{for } \rho < 0: U(x) = -x^\rho. \end{aligned} \tag{7.1}$$

832 It is well-known that the unit of payment is immaterial for this family. The most general
 833 family considered for $w(p)$ is Prelec's (1998) popular two-parameter family

⁸ We avoid the latter term because in nonexpected utility models as relevant for this paper, risk aversion depends not only on utility.

834 $w(p) = \left(\exp(-\beta(-\ln(p))^\alpha) \right). \quad (7.2)$

835 We will mostly use the one-parameter subfamily with $\beta=1$, as in Eq. 4.2.4, for reasons
 836 explained later. Substituting the above functions yields

837 $B(E) = \exp\left(-\left(\frac{-\ln(\frac{r(2r-r^2)^{1-p}}{(1-r)(1-r^2)^{1-p} + r(2r-r^2)^{1-p}})}{\beta}\right)^{1/\alpha}\right).$

838 for Eq. 4.3.6.

839 The model we estimate is as follows.

840 $R_{s,t,k}(j/20) = h(j/20, \alpha_t, \rho_t) + \varepsilon_{s,t,k}(j/20, \sigma_t^2). \quad (7.3)$

841 Here $R_{s,t,k}(j/20)$ is the reported probability of participant s for known probability $p=j/20$ ($10 \leq j \leq 19$) in treatment t ($t = \text{ALL}$ or $t = \text{ONE}$) for the k^{th} measurement for this probability, with
 842 only $k=1$ for $j = 10$, $k = 1,2$ for $11 \leq k \leq 18$, and $k = 1,2,3$ for $j = 19$. With β set equal to 1, α_t
 843 is the remaining probability-weighting parameter (Eq. 7.2), and ρ_t is the power of utility (Eq.
 844 7.1). The function h is the inverse of Eq. 4.2.5. Although we have no analytic expression for
 845 this inverse, we could calculate it numerically in the analyses. The error terms $\varepsilon_{s,t,k}(j/20)$ are
 846 drawn from a truncated normal distribution with mean 0 and treatment dependent variance
 847 σ_t^2 . The distribution of the error terms is truncated because reported probabilities below 0
 848 and above 1 are excluded by design. Error terms are identically and independently
 849 distributed across participants and choices. We employ maximum likelihood to estimate the
 850 parameters of Eq. 7.3. We also carried out an analysis at the individual level of the
 851 calibration part, with $\alpha_{s,t}$ and $\rho_{s,t}$ instead of α_t and ρ_t , i.e. with these parameters depending on
 852 the participant.

854 In the stock-price part, we tested for violations of additivity. With I the large interval of
 855 a stock, being the union $S \cup T$ of the two small intervals S and T , additivity of the uncorrected
 856 reported probabilities implies

857 $r_S + r_T = r_I. \quad (7.4)$

858 Hence, $r_S + r_T - r_I$ is an index of deviation from additivity, which we call the *additivity bias*
 859 of r . For the special case of S the universal event with r a decision-weighting function, Dow

860 & Werlang (1992) interpreted this quantitative index of nonadditivity as an index of
 861 uncertainty aversion.

862 Under the null hypothesis of additivity for risk-corrected reported probabilities B , binary
 863 additivity holds, and we can obtain $B(S) = 1 - B(S^c)$ for small intervals S in the experiment
 864 (cf. Eq. 2.2). Thus, under additivity of B , we have

865
$$B(S) + B(T) = B(I). \quad (7.5)$$

866 Hence, $B(S) + B(T) - B(I)$ is an index of deviation from additivity of B , and is B 's *additivity*
 867 *bias*.

868 We next discuss tests of the additivity bias. For each individual stock, and also for the
 869 average over all stocks, we tested for both treatments $t=ONE$ and $t=ALL$, (a) whether the
 870 additivity bias was zero or not, both with and without risk correction; (b) whether the
 871 additivity bias was enlarged or reduced by correction; (c) whether the absolute value of the
 872 additivity bias was enlarged or reduced by correction. We report only the tests for averages
 873 over all stocks.

874

875 **8. Results of the Calibration Part**

876 Risk-corrections and, in general, QSR measurements, do not make sense for participants who
 877 are hardly responsive to probabilities, so that $R(p)$ is almost flat on its entire domain. Hence
 878 we kept only those participants for whom the correlation between reported probability and
 879 objective probability is larger than 0.2 and, hence, dropped 4 participants. The following
 880 analyses are based on the remaining 89 participants.

881

882 *8.1. Group Averages*

883

884 We did several tests using Eq. 7.2 with β as a free (treatment-dependent or -independent)
 885 variable, but β 's estimates added little extra explanatory power to the other parameters and
 886 usually were close to 1. Hence, we chose to focus on a more parsimonious model in which
 887 the restriction $\beta_{ONE} = \beta_{ALL} = 1$ is employed. Table 8.1 lists the estimates for the model of Eq.
 888 7.3 for $\beta=1$ (Eq. 4.2.4 instead of Eq. 7.2) together with the estimates of some models with

889 additional restrictions. We first give results for group averages, assuming homogeneous
 890 participants.

891

892 TABLE 8.1. Estimation results at the aggregate level

Row	Restrictions	σ_{ONE}	α_{ONE}	ρ_{ONE}	σ_{ALL}	α_{ALL}	ρ_{ALL}	$-\text{LogL}$
1	NA	11.16* (0.30)	0.91* (0.06)	0.89* (0.14)	10.63* (0.30)	0.85* (0.04)	1.41* (0.07)	6513.84
2	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}}$ $= \rho_{\text{ONE}} = \rho_{\text{ALL}} = 1$	12.14* (0.31)	—	—	10.30* (0.26)	—	—	6554.55
3	$\alpha_{\text{ONE}} = \rho_{\text{ONE}} = 1$	12.14* (0.31)	—	—	10.63* (0.30)	0.85* (0.04)	1.41* (0.07)	6539.04
4	$\alpha_{\text{ALL}} = \rho_{\text{ALL}} = 1$	11.16* (0.30)	0.91* (0.06)	0.89* (0.14)	10.30* (0.27)	—	—	6529.36
5	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}}$	11.21* (0.30)	0.87* (0.03)	0.99* (0.08)	10.60* (0.29)	—	1.37* (0.06)	6514.31
6	$\rho_{\text{ONE}} = \rho_{\text{ALL}}$	11.40* (0.31)	0.79* (0.03)	1.19* (0.07)	10.47* (0.28)	0.96* (0.04)	—	6520.51
7	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}} = 1$	11.12* (0.29)	—	0.70* (0.04)	10.52* (0.29)	—	1.14* (0.03)	6519.68
8	$\rho_{\text{ONE}} = \rho_{\text{ALL}} = 1$	11.23* (0.29)	0.87* (0.02)	—	10.43* (0.28)	1.07* (0.02)	—	6522.46
9	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}} =$ $\rho_{\text{ONE}} = 1$	12.14* (0.31)	—	—	10.52* (0.29)	—	1.14* (0.03)	6544.09
10	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}} =$ $\rho_{\text{ALL}} = 1$	11.12* (0.29)	—	0.70* (0.04)	10.30* (0.27)	—	—	6530.14
11	$\alpha_{\text{ONE}} = \alpha_{\text{ALL}} = 1,$ $\rho_{\text{ONE}} = \rho_{\text{ALL}}$	12.05* (0.34)	—	0.98* (0.03)	10.30* (0.27)	—	—	6554.33

893 Standard errors in parentheses, * denotes significance at the 1% level.

894

895 *Overall need for risk-correction.* The 1st row of Table 8.1 shows the results for the most
 896 general model. The 2nd row presents the results without any correction. The likelihood
 897 reduces significantly (Likelihood Ratio test, $p = 0.01$) and substantially, so that risk-
 898 correction is called for. Risk-correction is also called for in both treatments in isolation, as
 899 the 3rd and 4th rows show, which significantly improve the likelihood relative to the 2nd row
 900 (Likelihood Ratio test; $p = 0.01$ for $t=\text{ALL}$, comparing 3rd to 2nd row; $p = 0.01$ for $t=\text{ONE}$,
 901 comparing 4th to 2nd row).

902

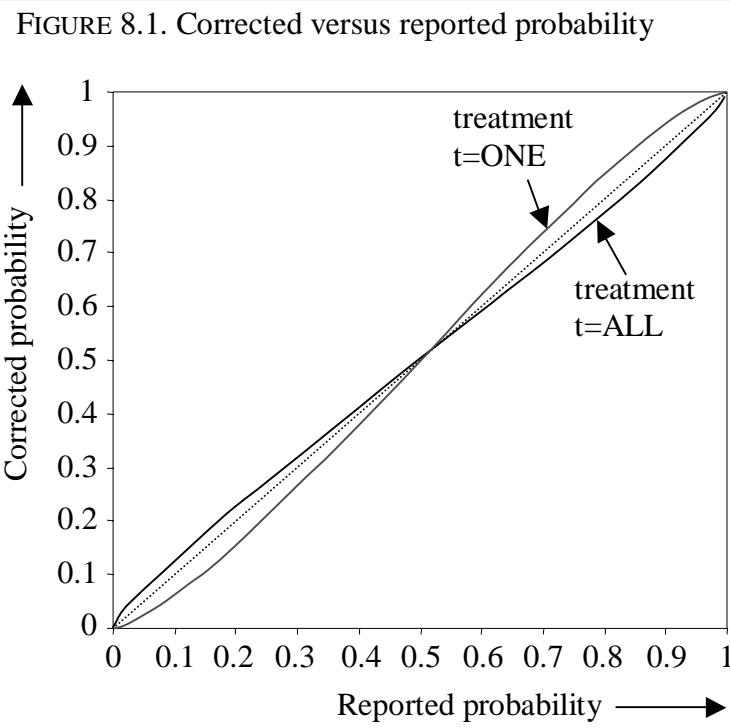
903 *Comparing the two treatments.* The likelihood for correcting only t=ALL (3rd row) is worse
 904 than for correcting only t=ONE (4th row), suggesting that there is more need for risk-
 905 correction for treatment t=ONE than for t=ALL. This difference does not seem to be caused
 906 by different probability weighting. The coefficients for probability weighting (α_{ONE} , α_{ALL}) in
 907 the 1st row are close to each other and are both smaller than 1. Apparently, probability
 908 weighting does not differ between t=ONE and t=ALL. Indeed, adding the restriction $\alpha_{\text{ONE}} =$
 909 α_{ALL} (5th row) does not decrease the likelihood of the data significantly (Likelihood ratio test;
 910 $p > 0.05$).

911 The difference between the two treatments is apparently caused by curvature of utility,
 912 captured by ρ_{ONE} and ρ_{ALL} . We obtain $\rho_{\text{ONE}} < \rho_{\text{ALL}}$: when only one decision is paid out then
 913 participants exhibit more concave curvature of utility than when all decisions are paid out.
 914 Given same probability weighting, it implies more risk aversion for t=ONE than for t=ALL
 915 (and R closer to 0.5). The finding is supported by comparing the 6th row of Table 8.1, with
 916 the restriction $\rho_{\text{ONE}} = \rho_{\text{ALL}}$, to the 1st row. This restriction significantly reduces the
 917 likelihood of observing the data (Likelihood Ratio test, $p = 0.01$).
 918

919 *Comparing utility and probability weighting.* Correcting only for utility curvature (7th row,
 920 $\alpha_{\text{ONE}} = \alpha_{\text{ALL}} = 1$) has a somewhat better likelihood than correcting only for probability
 921 weighting (8th row, $\rho_{\text{ONE}} = \rho_{\text{ALL}} = 1$).
 922

923 *Discussion of comparison of utility curvature and probability weighting for group-averages.*
 924 In deterministic choice, α could be determined through the flat part of R around 0.5, after
 925 which ρ could serve to improve the fit elsewhere. Statistically, however, α and ρ have much
 926 overlap, with risk aversion enhanced and R(p) moved towards 0.5 by increasing α and
 927 decreasing ρ , and one does not add much explanatory power to the other. It is, therefore,
 928 better to use only one of these parameters. Another reason to use only one concerns the
 929 individual analysis reported in the following subsection. Because we only have 20 choices
 930 per participant it is important to economize on the number of free parameters there.

931 We found above that ρ has a slightly better explanatory power than α . For this reason,
 932 and for reasons of convenience (see discussion section), we will only use the parameter ρ ,
 933 and assume $\alpha = 1$ henceforth. Figure 8.1 displays the resulting average risk-correction for
 934 the two treatments separately.
 935



935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952 Comparing the two treatments when there is no probability weighting. The average effect of
953 correction for utility curvature is not strong, especially for $t=ALL$. Yet this correction has a
954 significant effect, as can be seen from comparing the 7th row (general ρ) in Table 8.1 to its 9th
955 row ($\rho_{ALL} = 1$) (Likelihood Ratio test, $p = 0.01$).

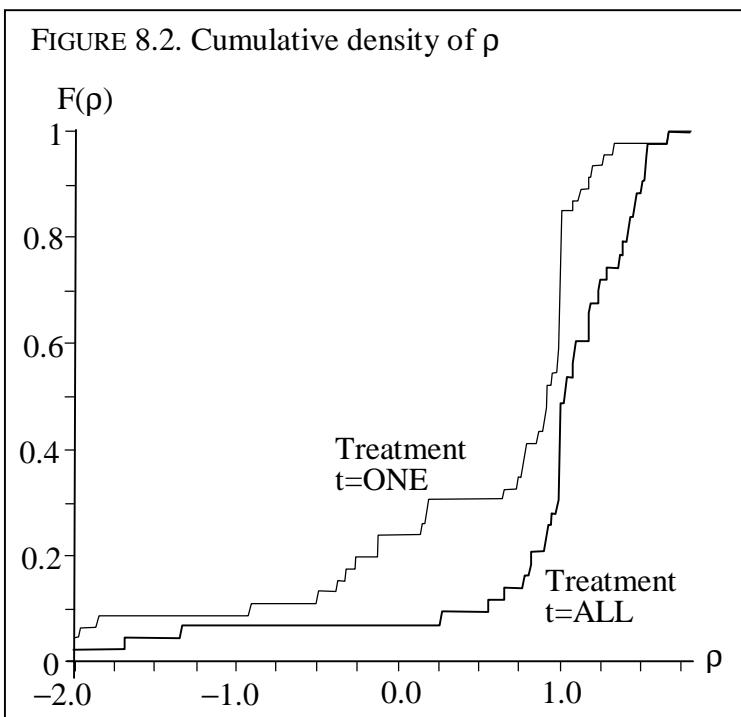
956

957

958

959 Need for risk-correction at the individual level. There is considerable heterogeneity in each
960 treatment. Whereas the corrections required were significant but small at the level of group
961 averages, they are big at the individual level. This appears from Figure 8.2, which displays
962 the cumulative distribution of the (per-subject) estimated ρ -coefficients for each treatment,
963 assuming $\alpha = \beta = 1$. There are wide deviations from the value $\rho=1$ (i.e., no correction) on
964 both sides. As seen from the group-average analysis, there are more deviations at the risk-
965 averse side of $\rho < 1$.

966



981 *Comparing the two treatments.* The ρ -coefficient distribution of treatment $t=ONE$ dominates
 982 the ρ -coefficient distribution of treatment $t=ALL$. A two-sided Mann-Whitney test rejects
 983 the null-hypothesis that the ranks of ρ -coefficients are equal across the treatments in favor of
 984 the hypothesis that the ρ -coefficients for $t=ONE$ are lower than for $t=ALL$ ($p=0.001$). It
 985 confirms that for group averages there is more risk aversion, moving R in the direction of 0.5,
 986 for $t=ONE$ than for $t=ALL$. The figure also shows that in an absolute sense there is more
 987 deviation from $\rho=1$ for $t=ONE$ than for $t=ALL$, implying that there are more deviations from
 988 expected value and more risk corrections for $t=ONE$ than for $t=ALL$.

989

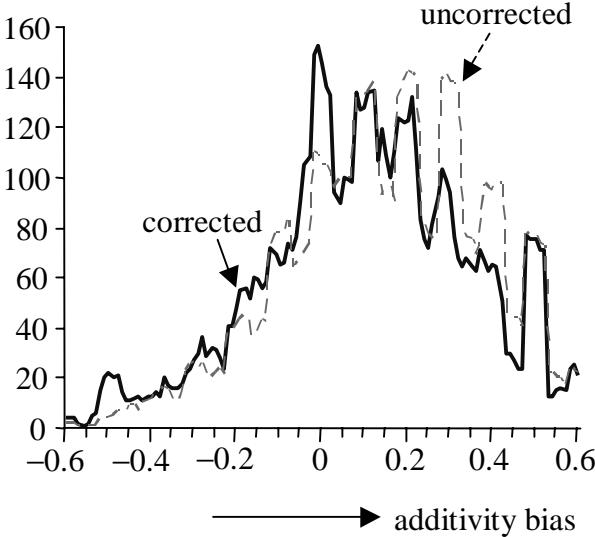
990 Unlike the median ρ -coefficients that are fairly close to each other for the two treatments
 991 (0.92 for $t=ONE$ versus 1.04 for $t=ALL$), the mean ρ -coefficients are substantially different
 992 (0.24 for $t=ONE$ versus 0.91 for $t=ALL$), which is caused by skewedness to the left for
 993 $t=ONE$. That is, there is a relatively high number of strongly risk-averse participants for
 994 $t=ONE$. Analyses of the individual ρ parameters (two-sided Wilcoxon signed rank sum tests)
 995 confirm findings of group-average analyses in the sense that the ρ -coefficients are
 996 significantly smaller than 1 for $t=ONE$ ($z = -3.50$, $p = 0.0005$), but not for $t=ALL$ ($z = 1.42$,
 997 $p = 0.16$).

998

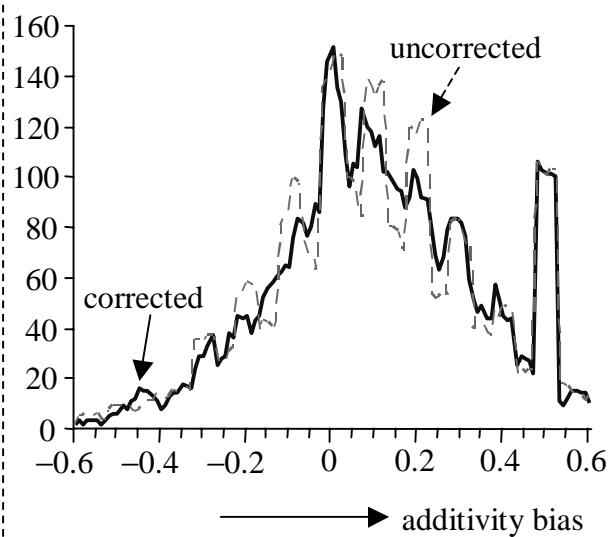
999 **9. Results for the Stock-Price Part: Risk-Correction and**
 1000 **Additivity**

1001
 1002 FIGURE 9.1. Empirical density of additivity bias for the two treatments

1003 FIG. a. Treatment t=ONE



1004 FIG. b. Treatment t=ALL



1005 For each interval $[\frac{j-2.5}{100}, \frac{j+2.5}{100}]$ of length 0.05 around $\frac{j}{100}$, we counted the number of
 1006 additivity biases in the interval, aggregated over 32 stocks and 89 individuals, for both
 1007 treatments. With risk-correction, there were 65 additivity biases between 0.375 and 0.425 in
 1008 the treatment t=ONE, and without risk-correction there were 95 such; etc.

1009 All comparisons hereafter are based on two-sided Wilcoxon signed rank sum tests.

1010 Figure 9.1 displays data, aggregated over both stocks and individuals, of the additivity biases
 1011 for t=ONE and for t=ALL. The figures show that the additivity bias is more often positive
 1012 than negative. Indeed, for virtually all stocks the additivity bias is significantly positive for
 1013 both treatments, showing in particular that additivity does not hold. This also holds when
 1014 taking the average additivity bias over all stocks as one data point per participant ($z = 5.27$, $p < 0.001$ for t=ONE, $z = 4.35$, $p < 0.001$ for t=ALL). We next consider whether risk
 1015 corrections reduce the violations of additivity.

1027 We first consider $t=ONE$. Here the risk corrections reduce the average additivity bias
 1028 significantly for 27 of the 32 stocks, and enlarge it for none. We only report the statistics for
 1029 the average additivity bias over all stocks per individual, which has overall averages 0.163
 1030 (uncorrected) and 0.120 (corrected), with the latter significantly smaller ($z = 3.21, p = 0.001$).
 1031 For assessing the degree of irrationality (additivity-violation) at the individual level, the
 1032 absolute values of the additivity bias are interesting. For $t=ONE$, Figure 9.1 suggests that
 1033 these are smaller after correction, because on average the corrected curve is closer to 0 on the
 1034 x-axis. These absolute values were significantly reduced for 9 stocks and enlarged for none.
 1035 Again, we only report the statistics for the average absolute value of the additivity bias over
 1036 all stocks per individual, which has overall averages 0.239 (uncorrected) and 0.228
 1037 (corrected), with the latter significantly smaller ($z = 2.26, p = 0.02$).

1038 For $t=ALL$, risk corrections did not significantly alter the average additivity bias. More
 1039 precisely, it gave a significant increase for 3 stocks and a significant decrease for 1 stock,
 1040 which, for 32 stocks, suggests no systematic effect. The latter was confirmed when we took
 1041 for each individual the average additivity bias over all stocks, with no significant differences
 1042 generated by correction (average 0.128 uncorrected and average 0.136 corrected; $z = -1.64, p$
 1043 = 0.1). Similar results hold for absolute values of additivity biases, which gave a significant
 1044 increase for 1 stock and a significant decrease for no stock, where taking for each individual
 1045 the average additivity bias over all stocks (average 0.237 uncorrected and average 0.239
 1046 corrected; $z = -0.36, p = 0.7$) also gave no significant difference.

1047 Classifications of individuals according to whether they exhibited more positive or more
 1048 negative additivity biases, and to whether risk corrections improved or worsened the
 1049 additivity bias more often, confirmed the patterns obtained above through stockwise
 1050 analyses, and will not be reported.

1051 Risk correction reduces the additivity bias for treatment $t=ONE$ to a level similar to that
 1052 observed for $t=ALL$ (averages 0.120 and 0.136). The overall pattern is that beliefs for
 1053 $t=ONE$ after correction, and for $t=ALL$ both before and after correction, exhibit a similar
 1054 degree of violation of additivity, which is clearly different from zero. The additivity bias is
 1055 not completely caused by nonlinear risk attitudes when participants report probabilities, but
 1056 has a genuine basis in beliefs.

1057

1058 **10. Discussion**1059 *10.1. Discussion of Methods*

1060

1061 We chose the evaluation date (June 1, 1991) sufficiently long ago to ensure that
 1062 participants would be unlikely to recognize the stocks or have private information about
 1063 them. In addition, no numbers were displayed on the vertical axis, making it extra hard for
 1064 participants to recognize specific stocks. We, thus, ensured that participants based their
 1065 probability judgments entirely on the prior information about past performance of the stocks
 1066 given by us. Given the large number of questions it is unlikely that participants noticed that
 1067 the graphs were presented more than once (three times) for each stock. Indeed, in informal
 1068 discussions after the experiment no participant showed awareness of this point.

1069 In some studies in the literature, the properness of scoring rules is explained to
 1070 participants by stating that it is in their best interest to state their true beliefs, either without
 1071 further explanation, or with the claim added that they will thus maximize their “expected”
 1072 money. A drawback of this explanation is that expected value maximization is empirically
 1073 violated, which is the central topic of this paper (Section 3). We, therefore, used an
 1074 alternative explanation that relates properness for unique events to observed frequencies of
 1075 repeated events (Appendix C).

1076 Besides the family of Prelec (1998), other parametric families for weighting functions
 1077 have been used in the literature, such as the family of Tversky & Kahneman (1992), and the
 1078 one of Goldstein & Einhorn (1987, Eqs. 22–24). We used Prelec’s family because it
 1079 performs equally well empirically as the other families but is analytically more tractable, for
 1080 example because its inverse can be defined easily (cf. Corollary 4.3.4). In addition, it is the
 1081 only one having an axiomatic foundation.

1082 In pragmatic applications of our method, more tractable families can be used to fit the
 1083 reported probabilities than the decision-theory-based curves that we used. For example, in
 1084 Figure 6.2 we also used quadratic regression to find the curve $p = a + br + cr^2$ that best fits
 1085 the data. The curve is virtually indistinguishable from the decision-theoretic curve. This
 1086 observation, together with Corollary 5.1 demonstrating that we only need the readily
 1087 observable reported probabilities and not the actual utility function or probability weighting
 1088 function to apply our method, shows that applications of our method are easy. The
 1089 theoretical part of this paper, and the decision-theory based curve-fitting that we adopted,

1090 served to prove that our method is in agreement with modern decision theories. If this thesis
 1091 is accepted, and the only goal is to obtain risk-corrected reported probabilities, then one may
 1092 choose the pragmatic shortcuts just described.

1093

1094 *10.2. Discussion of Main Results*

1095

1096 The significantly positive additivity bias that we found in all analyses shows that the
 1097 separate intervals together receive more weight than their union. This finding agrees with
 1098 other empirical findings in the literature, and it underlies the subadditivity of support theory
 1099 (Tversky & Koehler 1994).

1100 After some theoretical debates about the random-lottery incentive system (Holt 1986), as
 1101 in our treatment t=ONE, the system was tested empirically and found to be incentive-
 1102 compatible (Starmer & Sugden 1991). It is today the almost exclusively used incentive
 1103 system for measurements of individual preferences (Holt & Laury 2002). Unlike repeated
 1104 payments it avoids income effects such as Thaler & Johnson's (1990) house money effect,
 1105 and the drift towards expected value and linear utility. For the purpose of measuring
 1106 individual preference, the treatment t=ONE is, therefore, preferable. When the purpose is,
 1107 however, to derive subjective probabilities from proper scoring rules, and no risk-correction
 1108 is possible, then a drift towards expected value is actually an advantage, because uncorrected
 1109 proper scoring rules assume expected value. This point agrees with our findings, where less
 1110 risk-correction was required for the t=ALL treatment.

1111 For some applications group averages of probability estimates are most relevant, such as
 1112 when aggregating expert judgments or predicting group behavior. Then our statistical results
 1113 regarding “non-absolute” values of reported probabilities are most relevant. For the
 1114 assessment of rationality at the individual level, absolute values of the additivity biases are
 1115 most relevant.

1116

1117 *10.3. Discussion of Further Results*

1118

1119 The lack of extra explanatory power of parameter β in Eq. 7.2 should come as no
 1120 surprise because β and α behave similarly on $[0.5,1]$, increasing risk aversion there. They
 1121 mainly deviate from one another on $[0,0.5]$, where β continues to enhance risk aversion but α

1122 enhances the inverse-S shape that is mostly found empirically. The domain [0,0.5] is,
 1123 however, not relevant to our study (Observation 4.2.2).

1124 We found that the risk correction through the utility curvature parameter ρ fitted the data
 1125 somewhat better than the correction through the probability-weighting parameter α . This
 1126 finding may be interpreted as some descriptive support for expected utility. Another reason
 1127 that we used ρ and not α in our main analysis is that ρ , and utility curvature, are more well-
 1128 known in the economic literature than probability weighting, and are more analytically
 1129 tractable with R^{-1} defined everywhere. Although ρ indeed reflects the power of utility *if*
 1130 *expected utility is assumed*, we caution against unqualified interpretations here, as in any
 1131 study of risk aversion. The parameter ρ may also capture risk aversion generated by
 1132 probability weighting, and possibly by other factors.

1133

1134 *10.4. General Discussion*

1135

1136 Under proper scoring rules, beliefs are derived solely from decisions, and Eq. 2.1 is
 1137 taken purely as a decision problem, where the only goal of the agent is to optimize the
 1138 prospect received. Thus, this paper has analyzed proper scoring rules purely from the
 1139 decision-theoretic perspective supported with real incentives, and has corrected only for
 1140 biases resulting therefrom. Many studies have investigated direct judgments of belief
 1141 without real incentives, and then many other aspects play a role, leading for instance to the
 1142 often found overconfidence. Such introspective effects are beyond the scope of this paper.

1143 A drawback of our risk-correction procedure is that it requires individual measurements
 1144 of QSRs for given probabilities. If it is not possible to obtain individual measurements, then
 1145 it will be useful to use best-guess corrections, for instance through averages obtained from
 1146 individuals as similar as possible. Thus, at least, the systematic error for the group average to
 1147 risk attitude has been corrected for as good as is possible without requiring extra
 1148 measurements. In this respect the average curves in our Figure 8.1 are reassuring for existing
 1149 studies, because these curves suggest that only small corrections were called for regarding the
 1150 group averages in our context.

1151 Allen (1988) proposed to avoid biases of the QSR due to nonlinear utility by paying in
 1152 the probability of winning a prize instead of paying in money, and this procedure was
 1153 implemented by McKelvey & Page (1990). The procedure, however, only works if expected
 1154 utility holds, and there is much evidence against this assumption. Indeed, Selten, Sadrieh, &

1155 Abbink (1999) showed empirically that payment in probability does not generate the desired
 1156 risk neutral behavior.

1157 **11. Conclusion**

1158 Applications of proper scoring rules to measure subjective beliefs have so far been based
 1159 on the assumption of expected value maximization. However, many empirical deviations of
 1160 this assumption have been found. We have provided a method to correct for such deviations,
 1161 and have proved theoretically that our method provides such corrections under the modern
 1162 theories of nonexpected utility. These theories are empirically more realistic than expected
 1163 value maximization.

1164 We have demonstrated the feasibility and empirical tractability of our method in an
 1165 experiment, where we used it to investigate some properties of quadratic proper scoring rules
 1166 and beliefs. In a treatment with one big incentive for one randomly selected decision, we
 1167 found systematic distortions of (uncorrected) reported probabilities. No systematic
 1168 distortions were found in a treatment with many repeated decisions and repeated small
 1169 payments. When applying our correction procedure, both treatments give similar deviations
 1170 from additivity for the (corrected) beliefs. This finding suggests that subjective beliefs are
 1171 genuinely nonadditive. It means that expected value and expected utility are violated at the
 1172 level of beliefs, and that beliefs and ambiguity attitudes cannot be expressed in terms of
 1173 traditional probabilities, in agreement with Ellsberg's demonstrations. More general
 1174 nonadditive measures, such as used in nonexpected utility theories, are called for. For belief
 1175 elicitations where no risk correction can be implemented, repeated decisions with repeated
 1176 small payments are preferable to single large payments.

1177

1178 **Appendix A. Proofs and Technical Remarks**

1179 In Eqs. 4.2.1 and 4.2.2, probability p has a different decision weight when it yields the
 1180 best outcome of the prospect ($r > 0.5$) than when it yields the worst ($r < 0.5$). Similarly, in
 1181 Eqs. 4.3.1 and 4.3.2, E has a different decision weight when it yields the highest outcome ($r >$
 1182 0.5) than when it yields the lowest outcome ($r < 0.5$). Such a dependency of decision weights
 1183 on the ranking position of the outcome is called *rank-dependence* in the literature.

1184 Under rank-dependence, the sum of the decision weights in the evaluation of a prospect
 1185 are 1 even though $w(B(E))$ is not additive in E . This property is necessary for the functional
 1186 that evaluates prospects to satisfy natural conditions such as stochastic dominance, which
 1187 explains why theoretically sound nonexpected utility models could only be developed after
 1188 the discovery of rank dependence, a discovery that was made independently by Quiggin
 1189 (1982) for the special case of risk and by Schmeidler (1989, first version 1982) for the
 1190 general context of uncertainty.

1191 For qsr-prospects in Eq. 2.1, every choice $r < 0$ is inferior to $r = 0$, and $r > 1$ is inferior to
 1192 $r = 1$. The optimization problem does not change if we allow all real r , instead of $0 \leq r \leq 1$.
 1193 Hence, solutions $r = 0$ or $r = 1$ hereafter can be treated as interior solutions, and they satisfy
 1194 the first-order optimality conditions.

1195

1196 PROOF OF OBSERVATION 4.1.3. If $r = 0.5$ then the marginal utility ratio in Eq. 4.1.2 is 1, and
 1197 $p = 0.5$ follows. For the reversed implication, assume risk aversion. Then $r > 0.5$ is not
 1198 possible for $p = 0.5$ because then the marginal utility ratio in Eq. 4.1.2 would be at least 1 so
 1199 that the right-hand side of Eq. 4.1.2 would at most be 0.5, contradiction $r > 0.5$. Applying
 1200 this finding to E^c and using Eq. 2.2, $r < 0.5$ is not possible either, and $r = 0.5$ follows.

1201 Under strong risk seeking, r may differ from 0.5 for $p = 0.5$. For example, if $U(x) = e^{2.5x}$,
 1202 then $r = 0.14$ and $r = 0.86$ are optimal, and $r = 0.5$ is a local infimum, as calculations can
 1203 show. The same optimal values of r result under nonexpected utility with linear U , and with
 1204 $w(0.5) = 0.86$. Such large w -values also generate risk seeking.

1205

1206 PROOF OF THEOREM 4.3.2. We write π for the decision weight $W(E)$. For optimality of
 1207 interior solutions r , the first-order optimality condition for Eq. 4.3.1 is that
 1208 $\pi U'(a-b(1-r)^2)2b(1-r) - (1-\pi)U'(a-br^2)2br = 0$,
 1209 implying

$$1210 \quad \pi(1-r)U'(a-b(1-r)^2) = (1-\pi)rU'(a-br^2) \tag{A.1}$$

1211 or $\pi U'(a-b(1-r)^2) = r \times (\pi U'(a-b(1-r)^2) + (1-\pi)U'(a-br^2))$, and Eq. 4.3.5 follows.

1212 \square

1213

1214 PROOF OF COROLLARY 4.3.4. Let $r > 0.5$ be optimal, and write $\pi = W(E)$. Then Eq. A.1
 1215 implies

1216 $\pi \times ((1-r)U'(a-b(1-r)^2) + rU'(a-br^2)) = rU'(a-br^2)$, implying

$$1217 \quad \pi = \frac{r}{r + (1-r)\frac{U'(a-b(1-r)^2)}{U'(a-br^2)}} \quad (\text{A.2})$$

1218 Applying w^{-1} to both sides yields the theorem. \square

1219

1220 In measurements of belief one first observes r , and then derives $B(E)$ from it. Corollary
1221 4.3.4 gave an explicit expression. In general, it does not seem to be possible to write r as an
1222 explicit expression of $B(E)$ or, in the case of objective probabilities with $B(E) = p$, of the
1223 probability p .

1224

1225 PROOF OF COROLLARY 5.1. Theorem 4.3.2 implies that the right-hand side of Eq. 4.3.5 is r
1226 both as is, and with p substituted for $B(E)$. Because Eq. 4.3.5 is strictly increasing in
1227 $w(B(E))$, and w is strictly increasing too, $p = B(E)$ follows. \square

1228

1229 Appendix B. Models for Decision under Risk and Uncertainty

1230 For binary (two-outcome) prospects with both outcomes nonnegative, as considered in
1231 QSRs, Eqs. 4.3.1 and 4.3.2 have appeared many times in the literature. Early references
1232 include Allais (1953, Eq. 19.1), Edwards (1954 Figure 3), and Mosteller & Nogee (1951, p.
1233 398). The convenient feature that binary prospects suffice to identify utility U and the
1234 nonadditive $w \circ B = W$ was pointed out by Ghirardato & Marinacci (2001), Gonzalez & Wu
1235 (2003), Luce (1991, 2000), Miyamoto (1988), Pfanzagl (1959, p. 287 top and middle), and
1236 Wakker & Deneffe (1996, p. 1143 and pp.1144-1145). Luce & Narens (1985, Theorems 7.1
1237 & 7.2.2) and Ghirardato, Maccheroni, & Marinacci (2005) showed that the form in Eqs. 4.3.1
1238 and 4.3.2 is essentially the only one for binary prospects that allows for interval-scale utility
1239 of outcomes.

1240 The convenient feature that virtually all existing decision theories agree on the
1241 evaluation of binary prospects was pointed out by Miyamoto (1988), calling Eqs. 4.3.1 and
1242 4.3.2 generic utility, and Luce (1991), calling these equations binary rank-dependent utility.
1243 It was most clearly analyzed by Ghirardato & Marinacci (2001), who called the equations the
1244 biseparable model. These three works also axiomatized the model. The agreement for binary

prospects was also central in many works by Luce (e.g., Luce, 2000, Chapter 3) and in Gonzalez & Wu (2003). Only for more than two outcomes, the theories diverge (Mosteller & Nogee 1951 p. 398; Luce 2000, introductions to Chapters 3 and 5). We next describe some of the mentioned decision theories. Because we consider only nonnegative outcomes, losses play no role, and we describe prospect theory only for gains hereafter.

We begin with decision under risk, with known objective probabilities $P(E)$. Expected utility (von Neumann & Morgenstern, 1944) is the special case where w is the identity and $B(E) = P(E)$. Kahneman & Tversky's (1979) original prospect theory, Quiggin's (1982) rank-dependent utility, and Tversky & Kahneman's (1992) new prospect theory concern the special case of $B(E) = P(E)$, where w now can be nonlinear. The case $B(E) = P(E)$ also includes Gul's (1991) disappointment aversion theory.

We next consider the more general case where no objective probabilities need to be given for all events E . Expected utility (Savage 1954) is the special case where B is an additive, now “subjective,” probability and w is the identity. Choquet expected utility (Schmeidler 1989) and cumulative prospect theory (Tversky & Kahneman 1992) start from the general weighting function W , from which B obviously results as $w^{-1}(W)$, with w the probability weighting function for risk. The multiple priors model (Gilboa & Schmeidler 1989, Wald 1950) results with $W(E)$ the infimum value $P(E)$ over all priors P . Under Machina & Schmeidler's (1992) probabilistic sophistication, B is an additive probability measure.

1265

1266 Appendix C. Experimental Instructions

1267 [Instructions are translated from Dutch and concern the instructions of Treatment t=ONE
 1268 only]
 1269 This experiment is about statements of which you do not know whether they are true or not.
 1270 An example is the statement that snow did fall in Amsterdam in March 1861. You do not
 1271 know for sure whether this statement is true or not. We will ask you to indicate how likely it
 1272 is for you that such a statement is true, using probability judgments expressed in percentages.
 1273 Perhaps you will, for example, attach a probability of 30% to the statement that it snowed in
 1274 March 1861 in Amsterdam. We will then determine a score for you with the help of the
 1275 added table *on paper*.

1276 According to the table, for a probability judgment of 30% you get score 5100 if the
 1277 statement is true (snow did fall in Amsterdam in March 1861). You get score 9100 if the
 1278 statement is not true (snow did not fall in Amsterdam in March 1861). If you give a different
 1279 probability judgment, you get different scores, as shown in the table. For example, if you
 1280 give a probability judgment of 100%, your score is 10000 if the statement is true (snow did
 1281 fall), and 0 if the statement is not true (snow did not fall). We now like to check whether the
 1282 table with the scores is clear.

1283 [Practice questions using the table]

1284 Your answers were right. We will now explain some further features of the table. If you
 1285 are certain that the statement is true, then it is best for you to give the maximum probability
 1286 judgment of 100% because that gives the maximum score 10000 for a true statement. Every
 1287 other answer than surely yields a lower score. If you are certain that the statement is not true,
 1288 then it is similarly best to give the minimum probability judgment of 0%, because that gives
 1289 the maximum score 10000 for a false statement. In many cases you do not know for certain
 1290 whether a statement is true or not. We will now explain an important feature of the table on
 1291 the basis of a thought experiment.

1292

1293 **Thought experiment about repeated statements**

1294 The properties of the table can be well illustrated with the help of repeated statements.
 1295 Imagine, as a thought experiment, that you first have to give your probability judgment about
 1296 a particular statement (for example, snow in Amsterdam in a particular year, say 1861).
 1297 Imagine that you give judgment 30%, which means that you earn 5100 points in case of snow
 1298 and 9100 points in case of no snow. Next however, various repetitions of that statement are
 1299 being considered (snow in Amsterdam in March 1862, snow in Amsterdam in March 1863,
 1300, snow in Amsterdam in March 1960), leading to a total of 100 of such statements. For all
 1301 100 statements (thus every year between 1861 and 1960) your score will be determined
 1302 according to the table and your probability judgment (that is the same for every 100
 1303 statements). Your total score is then equal the sum of those 100 scores. For example, if it
 1304 did snow in Amsterdam in March 35 times in those 100 years, and it did not snow 65 times, a
 1305 probability judgment of 30% yields the following total score:

1306 $35 \times 5100 + 65 \times 9100 = 770000$

1307 We can also calculate this for other probability judgments, suppose that your probability
 1308 judgment was 35%, then your total score was:

1309 $35 \times 5775 + 65 \times 8775 = 772500$

1310 On the next page we show that your total-score is optimal if your probability judgment is
1311 exactly equal to that percentage. Put differently, if for example 35 of the 100 (35%)
1312 statements are true, then it is best for you to choose probability judgment 35% because it will
1313 give you the highest total-score.

1314

1315 **Now suppose that 35 of the 100 statements are true**

1316 We will determine what your total-score would have been at different judgments.

1317 [Table showing the total score for all possible probability judgments]

1318 It looks like judgment 35 is best. We conclude that if 35% of the statements are true,
1319 probability judgment 35 is optimal. Something similar holds for every percentage.

1320 **CONCLUSION.** For every percentage of true statements your total-score is optimal if you
1321 choose your probability judgment to be equal to that percentage. Check this for another
1322 number by clicking on continue.

1323

1324 **Now suppose that [entered number] of the 100 statements are true**

1325 We will determine what your total-score would have been at different judgments.

1326 [Table showing the total score for all possible probability judgments]

1327 It looks like judgment [entered number] is best. Thus we conclude that for [entered number]
1328 % true statements, probability judgment [entered number] is optimal. Something similar
1329 holds for every percentage.

1330 **CONFIRMATION OF THE CONCLUSION.** For every percentage of true statements
1331 your total-score is optimal if you choose your probability judgment to be equal to that
1332 percentage. If you want, you can check the conclusion again for another number than
1333 [entered number] by clicking on the link below.

1334

1335 **The experiment for non-repeated statements**

1336 The experiment we will perform concerns unique, and not repeated, statements. The various
1337 unique statements we consider are all different. For every single one of them you can give a
1338 different probability judgment.

1339 There is a big difference between the real experiment and the thought-experiment with
1340 repetition. In the thought experiment there was an objective-optimal probability judgment,
1341 based on the percentage of true statements. In the real experiment, there are no repetitions
1342 and for every probability judgment you get only one score.

1343 The thought experiment does give a guide for your probability judgment in the real
1344 experiment, with the percentage true statements as reference point. It is now based on your
1345 own subjective judgment however, and not on objective calculations. In the real experiment,
1346 there is no right or wrong answer. You purely choose what you like best.

1347 In the experiment, you will encounter all different sorts of statements, more or less
1348 probable ones, and you can choose all probability judgments ranging from 0% till 100%.
1349 You can only choose whole percentages.

1350

1351 **Payoff**

1352 This experiment consists of two parts. In both parts you will be asked to give probability
1353 judgments, 100 in part 1 and 20 in part 2. At the end of the experiment, one out of 120
1354 statements considered during the experiment will be randomly (with equal probability)
1355 selected and on the basis of your score at this statement you will be paid out in euros, where
1356 500 points is equal to 1 euro. Click on continue to read the instructions of the first part of the
1357 experiment.

1358

1359 **Instructions part 1**

1360 In the graph below you see the price of a stock from January till June in a year in the past.
1361 We used real stock prices of the Amsterdam Exchange when we made the graphs. The graph
1362 is scaled in such a way that the price of the stock always stays between the upper and lower
1363 axis. The same holds for the other graphs you will see later in this experiment. We consider
1364 the following statement: on the 31st of December in that particular year, the price of the stock
1365 in the graph was in the purple area. We ask you to give a probability judgment about the
1366 truth of this statement without any further information about the stock or the year. You can
1367 only base this on the course of the graph in the first half of the year.

1368 [Figure showing an example of graph of stock price]

1369 Your score at this question depends on your probability judgment and whether the statement
1370 is true or not, according to the table.

1371 [Figure showing the same graph but with three different end prices at 31st of December]

1372 The input of your probability judgment takes place in two phases: first you type in an integer
1373 number between 0 and 100, next you will be shown a menu in which your choice is
1374 reproduced with the corresponding scores from the table. At that moment you can still alter
1375 your choice and choose any other integer between 0 and 100. You can do this by selecting
1376 the up or down arrow, or by clicking the mouse in the menu and scroll to another probability

1377 judgment. Next, when you click on OK your choice is final and you continue with the next
 1378 statement. If you have any questions at this moment, raise your hand. The experimenter will
 1379 come to you.

1380

1381 **Instructions part 2**

1382 Part 1 of the experiment is now over. The second part of the experiment consists of 20
 1383 statements. Also in this part of the experiment you will be asked to give probability
 1384 judgments. The difference is that it does not concern the prediction of stock prices now, but
 1385 rolls with two 10-sided dice. On one of the dice are the values 00, 10, 20, 30, 40, 50, 60, 70,
 1386 80, 90 and on the other die are the values 1, 2, 3, 4, 5, 6, 7, 8, 9. Both dice will be rolled.
 1387 The sum of the outcomes has the values 1-100 (we consider the roll 00-0 as if it is 100),
 1388 where all values have the same probability.

1389 [Picture showing the two ten sided dice]

1390 An example of a statement is “the outcome is in the range 01-25.” This statement is true
 1391 when the outcome of the dice is indeed between 1 and 25 (including 25), and not true when
 1392 the outcome is higher than 25. The input of your probability judgment again takes place in
 1393 two phases: first you type in an integer number between 0 and 100, next you will be shown a
 1394 menu in which your choice is replicated with the corresponding scores from the table. At
 1395 that moment you can still alter your choice and choose any other integer number between 0
 1396 and 100. You can do this by selecting the up or down arrow, or by clicking the mouse in the
 1397 menu and scroll to another probability judgment. Next, when you click on OK your choice is
 1398 final and you continue with the next statement. Also in this part there is no right or wrong
 1399 answer; you again choose what you want best. At the end of the experiment one statement
 1400 will be selected and paid out. In case that this is a statement from part 2 of the experiment,
 1401 you will be asked to roll the two ten sided dice once.

1402 This is the end of part 2. Please raise your hand. The experimenter will come by so that
 1403 it can be determined which round will be paid out.

1404

1405 **References**

- 1406 Abdellaoui, Mohammed (2000), “Parameter-Free Elicitation of Utilities and Probability
 1407 Weighting Functions,” *Management Science* 46, 1497–1512.
 1408 Albers, Wulf, Robin Pope, Reinhard Selten, & Bodo Vogt (2000), “Experimental Evidence
 1409 for Attraction to Chance,” *German Economic Review* 1, 113–130.

- 1410 Allais, Maurice (1953), "Fondements d'une Théorie Positive des Choix Comportant un
 1411 Risque et Critique des Postulats et Axiomes de l'Ecole Américaine," *Colloques
 1412 Internationaux du Centre National de la Recherche Scientifique (Econométrie)* 40,
 1413 257–332. Paris: Centre National de la Recherche Scientifique. Translated into English,
 1414 with additions, as "The Foundations of a Positive Theory of Choice Involving Risk and a
 1415 Criticism of the Postulates and Axioms of the American School," in Maurice Allais &
 1416 Ole Hagen (1979, Eds), *Expected Utility Hypotheses and the Allais Paradox*, 27–145,
 1417 Reidel, Dordrecht, The Netherlands.
- 1418 Allen, Franklin (1987), "Discovering Personal Probabilities when Utility Functions are
 1419 Unknown," *Management Science* 33, 542–544.
- 1420 Becker, Gordon M., Morris H. de Groot, & Jacob Marschak (1964), "Measuring Utility by a
 1421 Single-Response Sequential Method," *Behavioral Science* 9, 226–232.
- 1422 Becker, Selwyn W. & Fred O. Brownson (1964), "What Price Ambiguity? Or the Role of
 1423 Ambiguity in Decision Making," *Journal of Political Economy* 72, 62–73.
- 1424 Bernoulli, Daniel (1738), "Specimen Theoriae Novae de Mensura Sortis," *Commentarii
 1425 Academiae Scientiarum Imperialis Petropolitanae* 5, 175–192. Translated into English
 1426 by Louise Sommer (1954), "Exposition of a New Theory on the Measurement of Risk,"
 1427 *Econometrica* 22, 23–36.
 1428 Reprinted in Alfred N. Page (1968, Ed.), *Utility Theory: A Book of Readings*, Chapter
 1429 11, Wiley, New York. Revised translation in William J. Baumol & Stephen M. Goldfeld
 1430 (Eds, 1968), *Precursors in Mathematical Economics: An Anthology*. Clowes and Sons,
 1431 London, Selection 2, 15–26.
- 1432 Bleichrodt, Han & José Luis Pinto (2000), "A Parameter-Free Elicitation of the Probability
 1433 Weighting Function in Medical Decision Analysis," *Management Science* 46,
 1434 1485–1496.
- 1435 Bliss, Robert R. & Nikolaos Panigirtzoglou (2004), "Option-Implied Risk Aversion
 1436 Estimates," *Journal of Finance* 59, 407–446.
- 1437 Braga, Jacinto & Chris Starmer (2005), "Preference Anomalies, Preference Elicitation, and
 1438 the Discovered Preference Hypothesis," *Environmental and Resource Economics* 32,
 1439 55–89.
- 1440 Brier, Glenn W. (1950), "Verification of Forecasts Expressed in Terms of Probability,"
 1441 *Monthly Weather Review* 78, 1–3.

- 1442 Camerer, Colin F. & Martin Weber (1992), "Recent Developments in Modelling Preferences:
 1443 Uncertainty and Ambiguity," *Journal of Risk and Uncertainty* 5, 325–370.
- 1444 Clemen, Robert T. & Kenneth C. Lichtendahl (2002), "Debiasing Expert Overconfidence: A
 1445 Bayesian Calibration Model, Fuqua School of Business, Duke University, Durham NC."
- 1446 Clemen, Robert T. & Fred Rolle (2001), "In Theory ... In Practice," *Decision Analysis*
 1447 *Newsletter* 20, No 1.
- 1448 de Finetti, Bruno (1931), "Sul Significato Soggettivo della Probabilità," *Fundamenta
 1449 Mathematicae* 17, 298–329. Translated into English as "On the Subjective Meaning of
 1450 Probability," in Paola Monari & Daniela Cocchi (Eds, 1993) "Probabilità e Induzione,"
 1451 Clueb, Bologna, 291–321.
- 1452 de Finetti, Bruno (1937), "La Prévision: Ses Lois Logiques, ses Sources Subjectives,"
 1453 *Annales de l'Institut Henri Poincaré* 7, 1–68. Translated into English by Henry E.
 1454 Kyburg Jr., "Foresight: Its Logical Laws, its Subjective Sources," in Henry E. Kyburg Jr.
 1455 & Howard E. Smokler (1964, Eds), *Studies in Subjective Probability*, 93–158, Wiley,
 1456 New York; 2nd edition 1980, 53–118, Krieger, New York.
- 1457 de Finetti, Bruno (1962), "Does It Make Sense to Speak of "Good Probability Appraisers""?
 1458 In Isidore J. Good (Ed.), *The Scientist Speculates: An Anthology of Partly-Baked Ideas*,
 1459 William Heinemann Ltd., London.
 1460 Reprinted as Chapter 3 in de Finetti, Bruno (1972), "Probability, Induction and
 1461 Statistics." Wiley, New York.
- 1462 Dempster, Arthur P. (1967), "Upper and Lower Probabilities Induced by a Multivalued
 1463 Mapping," *Annals of Mathematical Statistics* 38, 325–339.
- 1464 Dow, James & Sérgio R.C. Werlang (1992), "Uncertainty Aversion, Risk Aversion and the
 1465 Optimal Choice of Portfolio," *Econometrica* 60, 197–204.
- 1466 Echternacht, Gary J. (1972), "The Use of Confidence Testing in Objective Tests," *Review of
 1467 Educational Research* 42, 217–236.
- 1468 Edwards, Ward (1954), "The Theory of Decision Making," *Psychological Bulletin* 51,
 1469 380–417.
- 1470 Ellsberg, Daniel (1961), "Risk, Ambiguity and the Savage Axioms," *Quarterly Journal of
 1471 Economics* 75, 643–669.
- 1472 Fischer, Gregory W. (1982), "Scoring Rule Feedback and the Overconfidence Syndrome in
 1473 Subjective Probability Forecasting," *Organizational Behavior and Human Performance*
 1474 29, 357–369.

- 1475 Fischhoff, Baruch, Paul Slovic & Sarah Lichtenstein (1977), "Knowing with Certainty: The
 1476 Appropriateness of Extreme Confidence," *Journal of Experimental Psychology: Human
 1477 Perception and Performance* 3, 552–564.
- 1478 Ghirardato, Paolo, Fabio Maccheroni, & Massimo Marinacci (2005), "Certainty
 1479 Independence and the Separation of Utility and Beliefs," *Journal of Economic Theory*
 1480 120, 129–136.
- 1481 Ghirardato, Paolo & Massimo Marinacci (2001), "Risk, Ambiguity, and the Separation of
 1482 Utility and Beliefs," *Mathematics of Operations Research* 26, 864–890.
- 1483 Gilboa, Itzhak (1987), "Expected Utility with Purely Subjective Non-Additive Probabilities,"
 1484 *Journal of Mathematical Economics* 16, 65–88.
- 1485 Gilboa, Itzhak & David Schmeidler (1989), "Maxmin Expected Utility with a Non-Unique
 1486 Prior," *Journal of Mathematical Economics* 18, 141–153.
- 1487 Gilboa, Itzhak & David Schmeidler (1999), "A Theory of Case-Based Decisions." Cambridge
 1488 University Press, Cambridge, UK.
- 1489 Goldstein, William M. & Hillel J. Einhorn (1987), "Expression Theory and the Preference
 1490 Reversal Phenomena," *Psychological Review* 94, 236–254.
- 1491 Gonzalez, Richard & George Wu (1999), "On the Shape of the Probability Weighting
 1492 Function," *Cognitive Psychology* 38, 129–166.
- 1493 Gonzalez, Richard & George Wu (2003), "Composition Rules in Original and Cumulative
 1494 Prospect Theory," Graduate School of Business, University of Chicago, Chicago, IL.
- 1495 Good, Isidore J. (1952), "Rational Decisions," *Journal of the Royal Statistical Society Series*
 1496 *B* 14, 107–114.
- 1497 Guala, Francesco (2000), "Artefacts in Experimental Economics: Preference Reversals and
 1498 the Becker-DeGroot-Marschak Mechanism," *Economics and Philosophy* 16, 47–75.
- 1499 Gul, Faruk (1991), "A Theory of Disappointment Aversion," *Econometrica* 59, 667–686.
- 1500 Hanson, Robin (2002), Piece entitled "Wanna Bet?" in *Nature* 420, November 2002, pp.
 1501 354–355.
- 1502 Holt, Charles A. (1986), "Preference Reversals and the Independence Axiom," *American
 1503 Economic Review* 76, 508–513.
- 1504 Holt, Charles A. (2005), "Webgames and Strategy: Recipes for Interactive Learning." in
 1505 press.
- 1506 Holt, Charles A. & Susan K. Laury (2002), "Risk Aversion and Incentive Effects," *American
 1507 Economic Review* 92, 1644–1655.

- 1508 Huck, Steffen & Georg Weiszäcker (2002), "Do Players Correctly Estimate What Others Do?
 1509 Evidence of Conservatism in Beliefs," *Journal of Economic Behavior and Organization*
 1510 47, 71–85.
- 1511 Kadane, Joseph B. & Robert L. Winkler (1988), "Separating Probability Elicitation from
 1512 Utilities," *Journal of the American Statistical Association* 83, 357–363.
- 1513 Kahneman, Daniel & Amos Tversky (1979), "Prospect Theory: An Analysis of Decision
 1514 under Risk," *Econometrica* 47, 263–291.
- 1515 Karni, Edi & Zvi Safra (1987), "Preference Reversal and the Observability of Preferences by
 1516 Experimental Methods," *Econometrica* 55, 675–685.
- 1517 Keren, Gideon B. (1991), "Calibration and Probability Judgments: Conceptual and
 1518 Methodological Issues," *Acta Psychologica* 77, 217–273.
- 1519 Keynes, John Maynard (1921), "A Treatise on Probability." McMillan, London. Second
 1520 edition 1948.
- 1521 Knight, Frank H. (1921), "Risk, Uncertainty, and Profit." Houghton Mifflin, New York.
- 1522 Lehrer, Ehud (2001), "Any Inspection is Manipulable," *Econometrica* 69, 1333–1347.
- 1523 Liberman, Varda & Amos Tversky (1993), "On the Evaluation of Probability Judgments:
 1524 Calibration, Resolution, and Monotonicity," *Psychological Bulletin* 114, 162–173.
- 1525 Luce, R. Duncan (1991), "Rank- and-Sign Dependent Linear Utility Models for Binary
 1526 Gambles," *Journal of Economic Theory* 53, 75–100.
- 1527 Luce, R. Duncan (2000), "Utility of Gains and Losses: Measurement-Theoretical and
 1528 Experimental Approaches." Lawrence Erlbaum Publishers, London.
- 1529 Luce, R. Duncan & Louis Narens (1985), "Classification of Concatenation Measurement
 1530 Structures According to Scale Type," *Journal of Mathematical Psychology* 29, 1–72.
- 1531 Machina, Mark J. (1982), "Expected Utility" Analysis without the Independence Axiom,"
 1532 *Econometrica* 50, 277–323.
- 1533 Machina, Mark J. & David Schmeidler (1992), "A More Robust Definition of Subjective
 1534 Probability," *Econometrica* 60, 745–780.
- 1535 Manski, Charles F. (2004), "Measuring Expectations," *Econometrica* 72, 1329–1376.
- 1536 Marinacci, Massimo (2002), "Probabilistic Sophistication and Multiple Priors,"
 1537 *Econometrica* 70, 755–764.
- 1538 McClelland, Alastair & Fergus Bolger (1994), "The Calibration of Subjective Probabilities:
 1539 Theories and Models 1980–1994." In George Wright & Peter Ayton (Eds), *Subjective
 1540 Probability*, 453–481, Wiley, New York.

- 1541 McKelvey, Richard & Talbot Page (1990), "Public and Private Information: An
 1542 Experimental Study of Information Pooling," *Econometrica* 58, 1321–1339.
- 1543 Miyamoto, John M. (1988), "Generic Utility Theory: Measurement Foundations and
 1544 Applications in Multiatribute Utility Theory," *Journal of Mathematical Psychology* 32,
 1545 357–404.
- 1546 Mosteller, Frederick & Philip Nogee (1951), "An Experimental Measurement of Utility,"
 1547 *Journal of Political Economy* 59, 371–404.
- 1548 Murphy, Allan H. & Robert L. Winkler (1974), "Subjective Probability Forecasting
 1549 Experiments in Meteorology: Some Preliminary Results," *Bulletin of the American
 1550 Meteorological Society* 55, 1206–1216.
- 1551 Nyarko, Yaw & Andrew Schotter (2002), "An Experimental Study of Belief Learning Using
 1552 Elicited Beliefs," *Econometrica* 70, 971–1005.
- 1553 Palmer, Tim N. & Renate Hagedorn (2006, Eds), "Predictability of Weather and Climate."
 1554 Cambridge University Press, Cambridge.
- 1555 Pfanzagl, Johann (1959), "A General Theory of Measurement —Applications to Utility,"
 1556 *Naval Research Logistics Quarterly* 6, 283–294.
- 1557 Plott, Charles R. & Kathryn Zeiler (2005), "The Willingness to Pay-Willingness to Accept
 1558 Gap, the "Endowment Effect," Subject Misconceptions, and Experimental Procedures
 1559 for Eliciting Valuations," *American Economic Review* 95, 530–545.
- 1560 Prelec, Drazen (1998), "The Probability Weighting Function," *Econometrica* 66, 497–527.
 1561 (First version: Prelec, Drazen (1989), "On the Shape of the Decision Weight Function,"
 1562 Harvard Business School, Harvard University, Cambridge, MA, USA.)
- 1563 Prelec, Drazen (2004), "A Bayesian Truth Serum for Subjective Data," *Science* 306, October
 1564 2004, 462–466.
- 1565 Quiggin, John (1982), "A Theory of Anticipated Utility," *Journal of Economic Behaviour
 1566 and Organization* 3, 323–343.
- 1567 Raiffa, Howard (1968), "Decision Analysis." Addison-Wesley, London.
- 1568 Sandroni, Alvaro, Rann Smorodinsky, & Rakesh V. Vohra (2003), "Calibration with Many
 1569 Checking Rules," *Mathematics of Operations Research* 28, 141–153.
- 1570 Savage, Leonard J. (1954), "The Foundations of Statistics." Wiley, New York. (Second
 1571 edition 1972, Dover Publications, New York.)
- 1572 Savage, Leonard J. (1971), "Elicitation of Personal Probabilities and Expectations," *Journal
 1573 of the American Statistical Association* 66, 783–801.

- 1574 Schmeidler, David (1989), "Subjective Probability and Expected Utility without Additivity,"
 1575 *Econometrica* 57, 571–587.
- 1576 Seltén, Reinhard, Abdolkarim Sadrieh, & Klaus Abbink (1999), "Money Does not Induce
 1577 Risk Neutral Behavior, but Binary Lotteries Do even Worse," *Theory and Decision* 46,
 1578 211–249.
- 1579 Shafer, Glenn (1976), "A *Mathematical Theory of Evidence*." Princeton University Press,
 1580 Princeton NJ.
- 1581 Spiegelhalter, David J. (1986), "Probabilistic Prediction in Patient Management and Clinical
 1582 Trials," *Statistics in Medicine* 5, 421–433.
- 1583 Staël von Holstein, Carl-Axel S. (1972), "Probabilistic Forecasting: An Experiment Related
 1584 to the Stock Market," *Organizational Behaviour and Human Performance* 8, 139–158.
- 1585 Starmer, Chris & Robert Sugden (1991), "Does the Random-Lottery Incentive System Elicit
 1586 True Preferences? An Experimental Investigation," *American Economic Review* 81,
 1587 971–978.
- 1588 Thaler, Richard H. & Eric J. Johnson (1990), "Gambling with the House Money and Trying
 1589 to Break Even: The Effects of Prior Outcomes on Risky Choice," *Management Science*
 1590 36, 643–660.
- 1591 Tversky, Amos & Daniel Kahneman (1992), "Advances in Prospect Theory: Cumulative
 1592 Representation of Uncertainty," *Journal of Risk and Uncertainty* 5, 297–323.
- 1593 Tversky, Amos & Derek J. Koehler (1994), "Support Theory: A Nonextensional
 1594 Representation of Subjective Probability," *Psychological Review* 101, 547–567.
- 1595 von Neumann, John & Oskar Morgenstern (1944, 1947, 1953), "Theory of Games and
 1596 Economic Behavior." Princeton University Press, Princeton NJ.
- 1597 Wakker, Peter P. (2004), "On the Composition of Risk Preference and Belief," *Psychogical
 1598 Review* 111, 236–241.
- 1599 Wakker, Peter P. & Daniel Deneffe (1996), "Eliciting von Neumann-Morgenstern Utilities
 1600 when Probabilities Are Distorted or Unknown," *Management Science* 42, 1131–1150.
- 1601 Wald, Abraham (1950), "Statistical Decision Functions." Wiley, New York.
- 1602 Winkler, Robert L. (1967), "The Assessment of Prior Distributions in Bayesian Analysis,"
 1603 *Journal of the American Statistical Association* 62, 776–800.
- 1604 Winkler, Robert L. & Allan H. Murphy (1970), "Nonlinear Utility and the Probability
 1605 Score," *Journal of Applied Meteorology* 9, 143–148.

- 1606 Wright, William F. (1988), "Empirical Comparison of Subjective Probability Elicitation
1607 Methods," *Contemporary Accounting* 5, 47–57.
1608 Yates, J. Frank (1990), "*Judgment and Decision Making.*" Prentice Hall, London.
1609