

1 **A Truth-Serum for Non-Bayesians: Correcting Proper**
2 **Scoring Rules for Risk Attitudes¹**

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

4 a: CREED, Dept. of Economics, University of Amsterdam, Roetersstraat 11,

5 Amsterdam, 1018 WB, The Netherlands

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

7 Netherlands

8
9 October, 2007

10
11 ABSTRACT. Proper scoring rules provide convenient and highly efficient tools for eliciting
12 subjective beliefs. As traditionally used, however, they are valid only under expected value
13 maximization. This paper shows how proper scoring rules can be generalized to modern
14 (“nonexpected utility”) theories of risk and ambiguity, yielding mutual benefits: the empirical
15 realism of nonexpected utility is introduced in proper scoring rules, and the beauty and
16 efficiency of proper scoring rules is introduced in nonexpected utility. An experiment
17 demonstrates the feasibility of our generalized measurement instrument, yielding plausible
18 empirical results.

19
20 KEY WORDS: belief measurement, proper scoring rules, ambiguity, Knightian uncertainty,
21 subjective probability, nonexpected utility

22 JEL-CLASSIFICATION: D81, C60, C91

¹ A preliminary version of this paper circulated with the title “Is the Quadratic Scoring Rule Really Incentive Compatible?” This paper received helpful and detailed comments from Glenn Harrison.

23 **1. Introduction**

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

32 The measurement of subjective beliefs is important in many domains (Gilboa &
33 Schmeidler 1999; Machina & Schmeidler 1992; Manski 2004), and proper scoring rules have
34 been widely used accordingly, in accounting (Wright 1988), Bayesian statistics (Savage 1971),
35 business (Staël von Holstein 1972), education (Echternacht 1972), finance (Shiller, Kon-Ya, &
36 Tsutsui 1996), medicine (Spiegelhalter 1986), politics (Tetlock 2005), psychology
37 (McClelland & Bolger 1994), and other fields (Hanson 2002; Johnstone 2006; Prelec 2004).
38 Proper scoring rules are especially useful for giving experts incentives to exactly reveal their
39 degrees of belief. They are commonly used, for instance, to measure the degree of belief of
40 weather forecasters and to improve their calibration (Palmer & Hagedorn 2006; Yates 1990).
41 They have recently become popular in experimental economics and game theory. The
42 quadratic scoring rule is the most popular proper scoring rule today (McKelvey & Page 1990;
43 Nyarko & Schotter 2002; Palfrey & Wang 2007), and is the topic of this paper.

44 Proper scoring rules were introduced independently by Brier (1950), Good (1952, p. 112),
45 and de Finetti (1962). They have traditionally been based on the assumption of expected value
46 maximization, i.e. risk neutrality. All applications up to today that we are aware of have
47 maintained this assumption. Empirically, however, many deviations from expected value
48 maximization have been observed, and this may explain why proper scoring rules have not
49 been used by modern decision theorists so far. The first deviation was pointed out by Bernoulli
50 (1738), who noted that risk aversion prevails over expected value, so that, under expected
51 utility, utility has to be concave rather than linear. Second, Allais (1953) demonstrated, for
52 events with known probabilities, that people can be risk averse in ways that expected utility
53 cannot accommodate, so that more general decision theories are called for with other factors

54 besides utility curvature (Kahneman & Tversky 1979; Quiggin 1982; Tversky & Kahneman
55 1992). Third, Keynes (1921), Knight (1921), and Ellsberg (1961) demonstrated the importance
56 of ambiguity for events with unknown probabilities (“Knightian uncertainty”). Then
57 phenomena occur that are fundamentally different than those for known probabilities, which
58 adds to the descriptive failure of expected value. Gilboa (1987), Gilboa & Schmeidler (1989),
59 Hogarth & Einhorn (1990), Schmeidler (1989), and Tversky & Kahneman (1992) developed
60 decision theories that incorporate ambiguity. A recent empirical study of the new phenomena
61 is Halevy (2007), and typical illustrations of the new implications for economic theory are in
62 Hansen, Sargent, & Tallarini (1999) and Mukerji & Tallon (2001).

63 It is high time that proper scoring rules be updated from the expected-value model as
64 assumed in the 1950s, when proper scoring rules were introduced, to the current state of the art
65 in decision theory, where violations of expected value have been widely documented. Such an
66 update, provided by this paper, brings mutual benefits for practitioners of proper scoring rules
67 and for the study of risk and ambiguity. For practitioners of proper scoring rules we show
68 how to improve the empirical performance and validity of their measurement instrument. For
69 studies of risk and ambiguity we show how to benefit from the very efficient measurement
70 instrument provided by proper scoring rules. Regarding the first benefit, we bring bad news
71 when describing the many empirical deviations from expected value that distort classical
72 proper scoring rules, but good news when we give quantitative assessments of those distortions
73 and ways to correct for them. In the experiment in this paper we will find no systematic biases
74 for a repeated-payment treatment, so that no correction may be needed for group averages then.
75 This can be further good news for classical applications of proper scoring rules. Regarding the
76 second benefit, we show how subjective beliefs and ambiguity attitudes can easily be isolated
77 from risk attitude, using the incentive compatibility and efficiency of proper scoring rules.

78 Our correction technique can be interpreted as a new calibration method (Keren 1991;
79 Yates 1990) that does not need many repeated observations, unlike traditional calibration
80 techniques (Clemen & Lichtendahl 2005). An efficient aspect of our method is that we need
81 not elicit the entire risk attitudes of agents so as to correct for them. For instance, we need
82 not go through an entire measurement of the utility and probability weighting functions to
83 apply our correction. Instead, we can immediately infer the correction from a limited set of
84 readily observable data (the “correction curve”; see later).

85 We emphasize that the biases that we correct for need not concern mistakes on the part
86 of our subjects. Deviations from risk neutrality need not be irrational and, according to some,
87 even deviations from Bayesian beliefs need not be irrational, nor is the ambiguity aversion

88 that may be implied (Gilboa & Schmeidler 1989). Thus, learning and incentives need not
89 generate the required corrections. These corrections concern empirical deficiencies of the
90 model of expected value, i.e. mistakes on the part of researchers analyzing the data.

91 We illustrate the feasibility of our method through an experiment where we measure the
92 subjective beliefs of participants about the future performance of stocks after provision of
93 information about past performance. The empirical findings confirm the usefulness of our
94 method. We find that violations of additivity of subjective beliefs are reduced but not
95 eliminated by our corrections. Thus, the classical measurements will contain violations of
96 additivity that are partly due to the incorrect assumption of expected value, but partly are
97 genuine. Subjective beliefs genuinely violate additivity, and cannot be modeled through
98 additive subjective probabilities.

99 From the Bayesian perspective, violations of additivity are undesirable. Because we can
100 measure such violations, we can investigate which of several implementations of proper
101 scoring rules best approximate the Bayesian ideal. To illustrate this point, we compared two
102 experimental treatments: (1) only one single large decision is randomly selected and paid for
103 real; (2) every decision is paid, and subjects earn the sum of (moderate) payments. Because of
104 the law of large numbers one expects the results of treatment (2), with repeated small
105 payments, to stay closer to expected value and Bayesianism than those of treatment (1) will.
106 This was confirmed in our experiment, where smaller corrections were required for the
107 repeated payments than for the single payment.

108 The analysis of this paper consists of three parts. The first part (§§3-5) considers various
109 modern theories of risk and ambiguity, and derives implications for proper scoring rules from
110 these theories. This part is of interest to practitioners of proper scoring rules because it
111 shows what distortions affect these rules. It is of interest to decision theorists because it
112 shows a new field of application.

113 The second part of the paper, §§6-7, applies the revealed-preference reversal technique
114 to the results of the first part. That is, we do not assume theoretical models to derive
115 empirical predictions therefrom, but we assume empirical observations and derive the
116 theoretical models from those. §6 presents the main result of this paper, showing how
117 subjective beliefs can be derived from observed choices in an easy manner. §7 contains a
118 simple example illustrating such a derivation at the individual level. It shows in particular
119 that many decision-theoretic details, presented in the first part to justify our correction
120 procedures, need not be studied when applying our method empirically. Readers interested
121 only in applying our method empirically can skip most of §§3-6, reading only §3 up to

122 Theorem 3.1 and Corollary 6.4. For practitioners of proper scoring rules, the second part of
 123 this paper then shows how beliefs can be derived from observed proper scoring rules under
 124 more realistic descriptive theories. We introduce so-called risk-corrections to correct for
 125 distortions. For the study of subjective (possibly non-Bayesian) beliefs and ambiguity
 126 attitudes, the second part of this paper shows how proper scoring rules can be used to
 127 measure and analyze these concepts efficiently. An observed choice in a proper scoring rule
 128 gives as much information as an observed indifference in a binary choice while avoiding the
 129 empirical difficulties pertaining to the latter.

130 The third part of the paper, §§8-11, presents an experiment where we implement our
 131 correction method. We present some preliminary findings on nonadditive beliefs and on
 132 different implementations of real incentives. For brevity, detailed examinations of empirical
 133 implementations of our method, of the descriptive and normative properness of additive
 134 subjective beliefs, the effects of real incentives, and also of interpretations of beliefs and
 135 ambiguity attitudes, are left as topics for future study. Our contribution is to show how those
 136 concepts can be measured. The experiment of this paper only serves to demonstrate the
 137 feasibility of empirical implementation of our theoretical contribution.

138 §8 contains methodological details. §9 presents results regarding the biases that we
 139 correct for, and §10 presents some implications of the corrections of such biases.
 140 Discussions and conclusions are in §§11-12. Appendix A presents proofs and technical
 141 results, Appendix B surveys the implications of modern decision theories for our
 142 measurements, and Appendix C presents details of the experimental instructions.

143

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

145 Let E denote an event of which an agent is uncertain about whether or not it obtains,
 146 such as snow in Amsterdam in March 1932, whether a stock's value will decrease during the
 147 next half year, whether a ball randomly drawn from 20 numbered balls will have a number
 148 below 5, whether the 100th digit of π is 3, and so on. The degree of uncertainty of an agent
 149 about E will obviously depend on the information that the agent has about E . Some agents
 150 may even know with certainty about some of the events. Most events will, however, be
 151 uncertain. For most uncertain events, no objective probabilities of occurrence are known,

152 and our decisions have to be based on subjective assessments, consciously or not, of their
153 likelihood.

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

160 Several methods have been used in the literature to measure the subjective degree of
161 belief of an agent in an event E . Mostly these have been derived from: (a) binary
162 preferences, which only give inequalities or approximations; (b) binary indifferences, which
163 are hard to elicit, e.g. through the complex Becker-DeGroot-Marschak mechanism (Braga &
164 Starmer 2005; Karni & Safra 1987) or bisection (Abdellaoui, Vossman, & Weber (2005)); (c)
165 introspection, which is not revealed-preference based let alone incentive-compatible. Proper
166 scoring rules provide an efficient and operational manner for measuring subjective beliefs
167 that deliver what the above methods seek to do while avoiding the problems mentioned.

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

$$170 \quad (E: 1-(1-r)^2, 1-r^2), \tag{2.1}$$

171 is offered to the agent, where $0 \leq r \leq 1$ is a number that the agent can choose freely. The
172 number chosen is a function of E , sometimes denoted r_E , and is called the (*uncorrected*)
173 *reported probability* of E . The reasons for this term will be explained later. More general
174 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
175 restrict our attention to $a = b = 1$. No negative payments can occur, so that the agent never
176 loses money. It is obvious that if the agent is certain that E will obtain, then he will
177 maximize $1-(1-r)^2$, irrespective of $1-r^2$, and will choose $r=1$. Similarly, $r=0$ is chosen if E
178 will certainly not obtain. The choice of $r = 0.5$ gives a riskless prospect, yielding 0.75 with
179 certainty. Increasing r increases the payment under E but decreases it under E^c . Under the
180 event that happens, the QSR pays 1 minus the squared distance between the reported
181 probability of a clairvoyant (who assigns probability 1 to the event that happens) and the
182 reported probability of the agent (r under E , $1-r$ under E^c). The following symmetry between
183 E and E^c will be crucial in later theories.

184

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

189

190 Because of Observation 2.1, we have

$$191 \quad r_{E^c} = 1 - r_E. \quad (2.2)$$

192

193 3. Proper Scoring Rules and Subjective Expected Value

194 The first two parts of our analysis concern a theoretical analysis of proper scoring rules.
 195 This section considers the model commonly assumed for proper scoring rules, from their
 196 introduction in the 1950s up to today: *subjective expected value* maximization. It means,
 197 first, that the agent assigns a subjective probability p to each event E.² Second, the agent
 198 maximizes expected value with respect to probabilities.

199 For QSRs and an event E with (subjective) probability $P(E) = p$, subjective expected
 200 value implies that the agent maximizes

$$201 \quad p \times (1 - (1-r)^2) + (1-p) \times (1-r^2) = 1 - p(1-r)^2 - (1-p)r^2. \quad (3.1)$$

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

$$206 \quad R(1-p) = 1 - R(p). \quad (3.2)$$

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 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.

207 The following theorem demonstrates that the QSR is incentive compatible. The theorem
208 immediately follows from the first-order optimality condition $2p(1-r) - 2r(1-p) = 0$ in Eq.
209 3.1. Second-order optimality conditions are verified throughout this paper and will not be
210 mentioned in what follows.

211

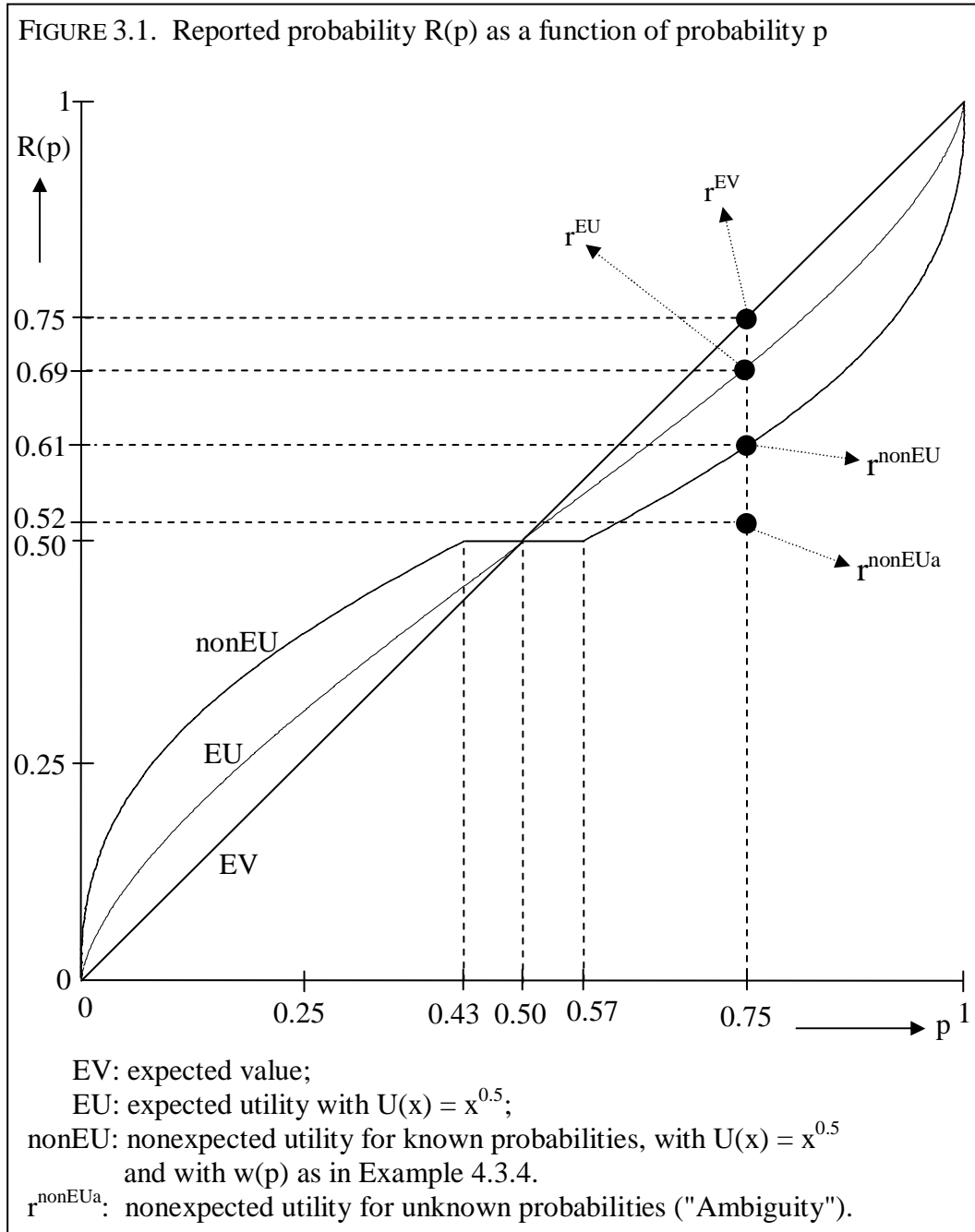
212 THEOREM 3.1. Under subjective expected value maximization, the optimal choice r_E is equal
213 to the probability p of event E , i.e. $R(p) = p$. \square

214

215 It is in the agent's best interest to truthfully report his subjective probability of E . This
216 explains the term “reported probability.” In Theorem 3.1, reported probabilities satisfy the
217 Bayesian additivity condition for probabilities. *Additivity* is the well-known property that the
218 probability of a disjoint union is the sum of the separate probabilities. We call the number r_E
219 the (*uncorrected*) *reported probability*.

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

224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257



EXAMPLE 3.2. An urn K ("known" distribution) contains 25 Crimson, 25 Green, 25 Silver,
 and 25 Yellow balls. One ball will be drawn at random. C designates the event of a crimson
 ball drawn, and G , S , and Y are similar. E is the event that the color is not crimson, i.e. it is
 the event $C^c = \{G, S, Y\}$. Under expected value maximization, $r_E = R(0.75) = 0.75$ is optimal
 in Eq. 2.1, yielding prospect $(E:0.9375, 0.4375)$ with expected value 0.8125. The point r_E is
 depicted as r^{EV} in Figure 3.1. Theorem 3.1 implies that $r_G = r_S = r_Y = 0.25$. We have $r_G + r_S$
 $+ r_Y = r_E$, and the reported probabilities satisfy additivity. \square

258 **4. Two Commonly Found Deviations from Expected Value under**
 259 **Risk, and Their Implications for Quadratic Proper Scoring Rules**

260 This section considers two factors that distort proper scoring rule measurements, and that
 261 should be corrected for. These factors concern decision attitudes and can be identified from
 262 decision under risk, with events for which probabilities are given. Proper scoring rules serve
 263 to examine other kinds of events, namely events with unknown probabilities. Those events
 264 will be the topic of the following sections. This section considers known probabilities only
 265 so as to identify biases, as a preparation for the following sections. §4.1 defines the domain
 266 of decision under risk, and then explains the organization of the other subsections.

267

268

4.1. Decision under Risk

269

270 ASSUMPTION 4.1.1. [Decision under Risk]. For event E , an objective probability p is given.

271 □

272

273 Any deviation of a reported probability r_E from the objective probability p entails a bias
 274 that should be corrected for. Under expected value maximization we obtain, similarly as in
 275 Theorem 3.1, that the agent should report $r=p$, so that there is no bias. The hypothetical
 276 situation of an agent using a subjective probability in Theorem 3.1 different than the
 277 objective probability in Assumption 4.1.1 cannot arise under plausible assumptions.³
 278 Subjective probabilities agree with objective probabilities whenever the latter exist, and this
 279 will be assumed throughout.

280 The effects of the factors that deviate from expected value and that distort the classical
 281 proper scoring rule measurements, explained later, are illustrated in Figure 3.1. Their

³ The first assumption is what defines decision under risk: that the only relevant aspect of events is their objective probability, and the second that we have sufficient richness of events to carry out the following reasoning. The claim then 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).

282 quantitative size will be illustrated through extensions of Example 3.2. §4.2 considers the
 283 first factor generating deviations, being nonlinear utility under expected utility. This section
 284 extends earlier studies of this factor by Winkler & Murphy (1970). We use expected utility
 285 and its primitives as in Savage's (1954) usual model. Extensions to alternative models
 286 (Broome 1990; Karni 2007; Luce 2000) are a topic for future research. §4.3 considers the
 287 second factor, namely violations of expected utility for known probabilities.

288

289 *4.2. The First Deviation: Utility Curvature*

290

291 Bernoulli (1738) put forward the first deviation from expected value. Because of the risk
 292 aversion in the so-called St. Petersburg paradox, Bernoulli proposed that people maximize
 293 the expectation of a *utility function* U . We assume that U is continuously differentiable with
 294 positive derivative everywhere, implying strict increasingness. We assume throughout that
 295 $U(0) = 0$. Eq. 3.1 is now generalized to

$$296 \quad pU(1-(1-r)^2) + (1-p)U(1-r^2) . \quad (4.2.1)$$

297 The first-order optimality condition for r , and a rearrangement of terms (as in the proof of
 298 Theorem 5.2), implies the following result. For $r \neq 0.5$, the theorem also follows as a
 299 corollary of Theorem 5.2 and Eq. 3.2.

300

301 THEOREM 4.2.1. Under expected utility with p the probability of event E , the optimal choice
 302 $r = R(p)$ satisfies:

$$303 \quad r = \frac{p}{p + (1-p)\frac{U'(1-r^2)}{U'(1-(1-r)^2)}} . \quad (4.2.2)$$

304 □

305

306 Figure 3.1 depicts an example of the function r under expected utility, indicated by the
 307 letters EU, and is similar to Figure 3 of Winkler & Murphy (1970); it is confirmed
 308 empirically by Huck & Weizsäcker (2002). The decision-based distortion in the direction of
 309 0.5 is opposite to the overconfidence (probability judgments too far from 0.5) mostly found
 310 in direct judgments of probability without real incentives (McClelland & Bolger 1994), and
 311 found among experts seeking to distinguish themselves (Keren 1991, p. 2f24 and 252; the

312 “expert bias”, Clemen & Rolle 2001). Optimistic and pessimistic distortions of probability
 313 can also result from nonlinear utility if the probability considered is a consensus probability
 314 for a group of individuals with heterogeneous beliefs (Jouini & Napp 2007).

315

316 EXAMPLE 4.2.2. Consider Example 3.2, but assume expected utility with $U(x) = x^{0.5}$.
 317 Substitution of Eq. 4.2.2 (or Theorem 5.2 below) shows that $r_E = R(0.75) = 0.69$ is optimal,
 318 depicted as r^{EU} in Figure 3.1, and yielding prospect (E:0.91, 0.52) with expected value
 319 0.8094. The extra risk aversion generated by concave U has led to a decrease of r_E by 0.06
 320 relative to Example 3.2, distorting the probability elicited, and generating an expected-value
 321 loss of $0.8125 - 0.8094 = 0.0031$. This amount can be interpreted as a risk premium,
 322 designating a profit margin for an insurance company. By Eq. 2.2, $r_C = 0.31$, and by
 323 symmetry $r_G = r_S = r_Y = 0.31$ too. The reported probabilities violate additivity, because $r_G +$
 324 $r_S + r_Y = 0.93 > 0.69 = r_E$. This violation in the data reveals that expected value does not
 325 hold. \square

326

327 OBSERVATION 4.2.3. Under expected utility with probability measure P , $r_E = 0.5$ implies
 328 $P(E) = 0.5$. Conversely, $P(E) = 0.5$ implies $r_E = 0.5$ if risk aversion holds. Under risk
 329 seeking, $r_E \neq 0.5$ is possible if $P(E) = 0.5$. \square

330

331 Theorem 4.2.1 clarifies the distortions generated by nonlinear utility, but it does not
 332 provide an explicit expression of $R(p)$, i.e. r as a function of p , or vice versa. It seems to be
 333 impossible, in general, to obtain an explicit expression of $R(p)$. We can, however, obtain an
 334 explicit expression of the inverse of $R(p)$, i.e. p in terms of r (Corollary 6.1). For numerical
 335 purposes, $R(p)$ can then be obtained as the inverse of that function—this is what we did in
 336 our numerical analyses, and how we drew Figure 3.1.

337

338 *4.3. The Second Deviation: Nonexpected Utility for Known Probabilities*

339

340 In the nonexpected utility analyses that follow, we will often restrict our attention to $r \geq$
 341 0.5. Results for $r < 0.5$ then follow by interchanging E and E^c , and the symmetry of
 342 Observation 2.1 and Eq. 2.2.

343 Event A is (*revealed*) *more likely than* event B if, for some positive outcome x , say $x =$
 344 100, the agent prefers $(A:x, 0)$ to $(B:x, 0)$. In all models considered hereafter, this
 345 observation is independent of the outcome $x > 0$. In view of the symmetry of QSRs in
 346 Observation 2.1, for $r \neq 0.5$ the agent will always allocate the highest payment to the most
 347 likely of E and E^c . It leads to the following restriction of QSRs.

348

349 OBSERVATION 4.3.1. Under the QSR in Eq. 2.1, the highest outcome is always associated
 350 with the most likely event of E and E^c . \square

351

352 Hence, QSRs do not give observations about most likely events when endowed with the
 353 worst outcome. Similar restrictions apply to all other proper scoring rules considered in the
 354 literature so far.

355 We now turn to the second deviation from expected value. With M denoting 10^6 , the
 356 preferences $M \succ (0.8: 5M, 0)$ and $(0.25:M, 0) \prec (0.20:5M, 0)$ are plausible. They would
 357 imply, under expected utility with $U(0) = 0$, the contradictory inequalities $U(M) > 0.8 \times$
 358 $U(5M)$ and $0.25U(M) < 0.20 \times U(5M)$ (implying $U(M) < 0.8 \times U(5M)$), so that they falsify
 359 expected utility. It has since been shown that this paradox does not concern an exceptional
 360 phenomenon pertaining only to hypothetical laboratory choices with extreme amounts of
 361 money, but that the phenomenon is relevant to real decisions for realistic stakes (Kahneman
 362 & Tversky 1979). The Allais paradox and other violations of expected utility have led to
 363 several alternative models for decision under risk, the so-called nonexpected utility models
 364 (Machina 1987; Starmer 2000; Sugden 2004). For the prospects relevant to this paper, QSRs
 365 with only two outcomes and no losses, all presently popular static nonexpected-utility
 366 evaluations of qsr-prospects (Eq. 2.1) are of the following form (see Appendix B). We first
 367 present such evaluations for the case of highest payment under event E , i.e. $r \geq 0.5$, which can
 368 be combined with $p \geq 0.5$.

$$369 \quad \text{For } r \geq 0.5: w(p)U(1-(1-r)^2) + (1-w(p))U(1-r^2). \quad (4.3.1)$$

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

$$374 \quad \text{For } r < 0.5: w(1-p)U(1-r^2) + (1-w(1-p))U(1-(1-r)^2). \quad (4.3.2)$$

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

377 Both in Eq. 4.3.1 and in Eq. 4.3.2, w is applied only to probabilities $p \geq 0.5$, and needs to
 378 be assessed only on this domain in what follows. This restriction is caused by Observation
 379 4.3.1. We display the implication.

380

381 OBSERVATION 4.3.2. For the QSR, only the restriction of w to $[0.5,1]$ plays a role, and w 's
 382 behavior on $[0,0.5)$ is irrelevant. \square

383

384 Hence, for the risk-correction introduced later, we need to estimate w only on $[0.5,1]$. An
 385 advantage of this point is that the empirical findings about w are uncontroversial on this
 386 domain, the general finding being that w underweights probabilities there.⁴

387

388 THEOREM 4.3.3. Under nonexpected utility with p the probability of event E , the optimal
 389 choice $r = R(p)$ satisfies:

$$390 \quad \text{For } r > 0.5: r = \frac{w(p)}{w(p) + (1-w(p)) \frac{U'(1-r^2)}{U'(1-(1-r)^2)}}. \quad (4.3.3)$$

391 \square

392

393 The above result, again, follows from the first-order optimality condition, and also
 394 follows as a corollary of Theorem 5.2 below. As an aside, the theorem shows that QSRs
 395 provide an efficient manner for measuring probability weighting on $(0.5, 1]$ if utility is linear,
 396 because then simply $r = R(p) = w(p)$. An extension to $[0, 0.5]$ can be obtained by a
 397 modification of QSRs, discussed further in the next section (Eqs. 5.6 and 5.7).

398

399 EXAMPLE 4.3.4. Consider Example 4.2.2, but assume nonexpected utility with $U(x) = x^{0.5}$
 400 and

$$401 \quad w(p) = \left(\exp(-(-\ln(p))^\alpha) \right) \quad (4.3.4)$$

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

402 with parameter $\alpha = 0.65$ (Prelec 1998). This function agrees with common empirical
 403 findings (Tversky & Kahneman 1992; Abdellaoui 2000; Bleichrodt & Pinto 2000; Gonzalez
 404 & Wu 1999). From Theorem 4.3.3 it follows that $r_E = R(0.75) = 0.61$ is now optimal,
 405 depicted as r^{nonEU} in Figure 3.1. It yields prospect (E:0.85, 0.63) with expected value 0.7920.
 406 The extra risk aversion relative to Example 4.2.2 generated by w for this event E has led to an
 407 extra distortion of r_E by 0.08. The extra expected-value loss (and, hence, the extra risk
 408 premium) relative to Example 4.2.2 is $0.8094 - 0.7920 = 0.0174$. By Eq. 4.3.1, $r_C = 0.39$,
 409 and by symmetry $r_G = r_S = r_Y = 0.39$ too. The reported probabilities strongly violate
 410 additivity, because $r_G + r_S + r_Y = 1.17 > 0.61 = r_E$. \square

411

412 Figure 3.1 illustrates the effects through the curve indicated by nonEU. The curve is flat
 413 around $p = 0.5$, more precisely, on the probability interval $[0.43, 0.57]$. For probabilities
 414 from this interval the risk aversion generated by nonexpected utility is so strong that the
 415 agent goes for maximal safety and chooses $r = 0.5$, corresponding with the sure outcome 0.75
 416 (cf. Manski 2004 footnote 10). Such a degree of risk aversion is not possible under expected
 417 utility, where $r = 0.5$ can happen only for $p = 0.5$ (Observation 4.2.3). This observation
 418 cautions against assigning specific levels of belief to observations $r = 0.5$, because proper
 419 scoring rules may be insensitive to small changes in the neighborhood of $p = 0.5$.

420 Up to this point, we have considered deviations from expected value and Bayesianism at
 421 the level of decision attitude, and beliefs themselves were not yet affected. This will be
 422 different in the next section.

423

424 **5. A Third Commonly Found Deviation from Subjective Expected** 425 **Value Resulting from Non-Bayesian Beliefs and Ambiguity, and** 426 **Its Implications for Quadratic Proper Scoring Rules**

427 This section considers a third deviation from expected value maximization. This
 428 deviation does not (merely) concern decision attitudes as did the two deviations examined in
 429 the preceding section. It rather concerns subjective beliefs for events with unknown
 430 probabilities (which involves ambiguity). These are the events that proper scoring rules serve
 431 to examine. Thus, the deviation in this section does not concern something we necessarily

432 have to correct for, but rather it concerns something that we want to measure and investigate
 433 without a commitment as to what it should look like.

434 In applications of proper scoring rules it is commonly assumed that the agent chooses
 435 (Bayesian) subjective probabilities $p = P(E)$ for such events, where these subjective
 436 probabilities are assumed to satisfy the laws of probability. The agent evaluates prospects the
 437 same way for subjective probabilities as if these probabilities were objective, leading to the
 438 following modification of Eq. 4.3.1:

$$439 \quad \text{For } r \geq 0.5: w(P(E))U(1-(1-r)^2) + (1-w(P(E)))U(1-r^2). \quad (5.1)$$

440 For w the identity with $w(P(E)) = P(E)$, Eq. 5.1 reduces to subjective expected utility, the
 441 subjective version of Eq. 4.2.1. In applications of proper scoring rules it is commonly
 442 assumed that not only w , but also U is the identity, leading to subjective expected value
 443 maximization, the model analyzed in §3.

444 The approach to unknown probabilities of Eq. 5.1, treating uncertainty as much as
 445 possible in the same way as risk, is called *probabilistic sophistication* (Machina &
 446 Schmeidler 1992). All results of §4 can be applied to this case, with distortions generated by
 447 nonlinear U and w . Probabilistic sophistication can be interpreted as a last attempt to
 448 maintain Bayesianism at least at the level of beliefs. Empirical findings, initiated by Ellsberg
 449 (1961), have demonstrated however that probabilistic sophistication is commonly violated
 450 empirically.

451

452 EXAMPLE 5.1 [Violation of Probabilistic Sophistication]. Consider Example 4.3.4, but now
 453 there is an additional urn A (“ambiguous”). Like urn K, A contains 100 balls colored
 454 Crimson, Green, Silver, or Yellow, but now the proportions of balls with these colors are
 455 unknown. C_a designates the event of a crimson ball drawn from A, and G_a , S_a , and Y_a are
 456 similar. E_a is the event $C_a^c = \{G_a, S_a, Y_a\}$. If probabilities are assigned to drawings from the
 457 urn A (as assumed by probabilistic sophistication) then, in view of symmetry, we must have
 458 $P(C_a) = P(G_a) = P(S_a) = P(Y_a)$, so that these probabilities must be 0.25. Then $P(E_a)$ must be
 459 0.75, as was $P(E)$ in Example 4.3.4. Under probabilistic sophistication combined with
 460 nonexpected utility as in Example 4.3.4, r_{E_a} must be the same as r_E in Example 4.3.4 for the
 461 known urn, i.e. $r_{E_a} = 0.61$. It implies that people must be indifferent between $(E:x, y)$ and
 462 $(E_a:x, y)$ for all x and y . The latter condition is typically violated empirically. People usually
 463 have a strict preference for known probabilities, i.e.

464 $(E:x, y) > (E_a:x, y)$.⁵

465 Consequently, it is impossible to model beliefs about uncertain events E_a through
 466 probabilities, and probabilistic sophistication fails. This observation also suggests that r_{E_a}
 467 may differ from r_E . \square

468
 469 The deviations from expected value revealed by Ellsberg through the above example
 470 cannot be explained by utility curvature or probability weighting, and must be generated by
 471 other factors. Those other, new, factors refer to components of beliefs and decision attitudes
 472 that are typical of unknown probabilities. They force us to give up on the additive measure
 473 $P(E)$ in our model. Besides decisions, also beliefs may deviate from the Bayesian principles.
 474 The important difference between known and unknown probabilities was first emphasized by
 475 Keynes (1921) and Knight (1921).

476 As explained in Appendix B, virtually all presently existing models for decision under
 477 uncertainty evaluate the qsr-prospect of Eq. 2.1 in the following way:

478 For $r \geq 0.5$: $W(E)U(1-(1-r)^2) + (1-W(E))U(1-r^2)$. (5.2)

479 Here W is a nonadditive set function often called *weighting function* or capacity, which
 480 satisfies the natural requirements that W assigns value 0 to the vacuous event \emptyset , value 1 to
 481 the universal event, and is increasing in the sense that $C \supset D$ implies $W(C) \geq W(D)$. For
 482 completeness, we also give the formula for $r < 0.5$, which can be obtained from Eq. 5.2
 483 through symmetry (Observation 2.1).

484 For $r < 0.5$: $(1-W(E^c))U(1-(1-r)^2) + W(E^c)U(1-r^2)$. (5.3)

485 Under probabilistic sophistication (Eq. 5.1), subjective belief P can be recovered from W
 486 through $P = w^{-1}(W)$, where w^{-1} can be interpreted as a correction for non-neutral risk
 487 attitudes. Machina (2004) argued that almost-objective probabilities can be constructed in
 488 virtually all circumstances of uncertainty, so that a domain for w is always available. In
 489 general, the “risk-corrected” function $w^{-1}(W)$ need not be a probability. We write

490 $B(E) = w^{-1}(W)$.

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.

491 B is what remains if the risk component w is taken out from W . It is common in
 492 decision theory to interpret factors beyond risk attitude as ambiguity. Then B reflects
 493 ambiguity attitude. There is no consensus about the extent to which ambiguity reflects non-
 494 Bayesian beliefs, and to what extent it reflects non-Bayesian decision attitudes beyond belief.
 495 If the equality $B(E) + B(E^c) = 1$ (*binary additivity*) is violated, then it can further be debated
 496 whether $B(E)$ or $1 - B(E^c)$ is to be taken as an index of belief or of ambiguity. Such
 497 interpretations have not yet been settled, and further studies are called for. We will usually
 498 refer to B as reflecting beliefs, to stay close to the terminology used in the literature on proper
 499 scoring rules today. On some occasions we will refer to the decision-theoretic ambiguity.
 500 Irrespective of the interpretation of B , it is clear that the behavioral component w^{-1} of risk
 501 attitude should be filtered out before an interpretation of belief can be considered. We show
 502 how this filtering out can be done.

503 In Schmeidler (1989), the main paper to initiate Eqs. 5.2 and 5.3, w was assumed linear,
 504 with expected utility for given probabilities, and W coincided with B . Schmeidler interpreted
 505 this component as reflecting beliefs. So did the first paper on nonadditive measures for
 506 decisions, Shackle (1949). Many studies of direct judgments of belief have supported the
 507 thesis that subjective beliefs may deviate from Bayesian probabilities (McClelland & Bolger
 508 1994; Shafer 1976; Tversky & Koehler 1994). Bounded rationality is an extra reason to
 509 expect violations of additivity at the level of beliefs (Aragones et al. 2005; Charness & Levin
 510 2005).

511 We rewrite Eq. 5.2 as

$$512 \quad \text{For } r \geq 0.5: w(B(E))U(1-(1-r)^2) + (1-w(B(E)))U(1-r^2). \quad (5.4).$$

$$513 \quad \text{For } r < 0.5: (1-w(B(E^c)))U(1-(1-r)^2) + w(B(E^c))U(1-r^2). \quad (5.5)$$

514 In general, B assigns value 0 to the vacuous event \emptyset , value 1 to the universal event, and B is
 515 increasing in the sense that $C \supset D$ implies $B(C) \geq B(D)$. These properties similarly hold for
 516 the composition $w(B(\cdot))$, as we saw above.

517 As with the weighting function w under risk, B is also applied only to the most likely one
 518 of E and E^c in the above equations, reflecting again the restriction of the QSR of Observation
 519 4.3.1. Hence, under traditional QSR measurements we cannot test binary additivity directly
 520 because we measure $B(E)$ only when E is more likely than E^c . These problems can easily be
 521 amended by modifications of the QSR. For instance, we can consider prospects

$$522 \quad (E: 2-(1-r)^2, 1-r^2), \quad (5.6)$$

523 i.e. qsr-prospects as in Eq. 2.1 but with a unit payment added under event E. The classical
 524 proper-scoring-rule properties of §2 are not affected by this modification, and the results of
 525 §3 are easily adapted. With this modification, we have the liberty to combine event E with
 526 the highest outcome both if E is more likely than E^c and if E is less likely, and we avoid the
 527 restriction of Observation 4.3.1. We then can observe w of the preceding subsection, and
 528 $W(E)$ and $B(E)$ over their entire domain. Similarly, with prospects

$$529 \quad (E: 1-(1-r)^2, 2-r^2), \quad (5.7)$$

530 we can measure the duals $1 - W(E^c)$, $1 - w(1-p)$, and $1 - B(E^c)$ over their entire domain. In
 531 this study we confine our attention to the QSRs of Eq. 2.1 as they are classically applied
 532 throughout the literature. We reveal their biases according to the current state of the art of
 533 decision theory, suggest remedies whenever possible, and signal the problems that remain.
 534 Further investigations of the, we think promising, modifications of QSRs in the above
 535 equations are left to future studies.

536 The restrictions of the classical QSRs will also hold for the experiment reported later in
 537 this paper. There an application of the QSR to events less likely than their complements are
 538 to be interpreted formally as the measurement of $1 - B(I^c)$. The restrictions also explain why
 539 the theorems below concern only the case of $r > 0.5$ (with $r = 0.5$ as a boundary solution).

540 The following theorem, our main theorem, specifies the first-order optimality condition
 541 for interior solutions of r for general decision making, incorporating all deviations described
 542 so far.

543

544 THEOREM 5.2. Under Eq. 5.4, the optimal choice r satisfies:

$$545 \quad \text{If } r > 0.5, \text{ then } r = r_E = \frac{w(B(E))}{w(B(E)) + (1-w(B(E))) \frac{U'(1-r^2)}{U'(1-(1-r)^2)}}. \quad (5.8)$$

546 □

547

548 We cannot draw graphs as in Figure 3.1 for unknown probabilities, because the x-axis
 549 now concerns events and not numbers. The W values of ambiguous events will be relatively
 550 low for an agent with a general aversion to ambiguity, so that the reported probabilities r in
 551 Eq. 5.8 will be relatively small, i.e. close to 0.5. We give a numerical example.

552

553 EXAMPLE 5.3. Consider Example 5.1. Commonly found preferences $(E:100, 0) \succ (E_a:100, 0)$
 554 imply that $w(B(E_a)) < w(B(E)) = w(0.75)$. Hence, by Theorem 5.2, r_{E_a} will be smaller than
 555 r_E . Given the strong aversion to unknown probabilities that is often found empirically
 556 (Camerer & Weber 1992), we will assume that $r_{E_a} = 0.52$. It is depicted as r^{nonEU_a} in Figure
 557 3.1, and yields prospect $(E_a:0.77, 0.73)$ with expected value 0.7596. The extra preference for
 558 certainty relative to Example 4.3.4 generated by unknown probabilities for this event E_a has
 559 led to an extra distortion of r_{E_a} by $0.61 - 0.52 = 0.09$. The extra expected-value loss relative
 560 to Example 4.3.4 is $0.7920 - 0.7596 = 0.0324$. This amount can be interpreted as the
 561 ambiguity-premium component of the total uncertainty premium. By Eq. 4.3.1, $r_C = 0.48$,
 562 and by symmetry $r_G = r_S = r_Y = 0.48$ too. The reported probabilities violate additivity to an
 563 extreme degree, because $r_G + r_S + r_Y = 1.44 > 0.52 = r_{E_a}$. The behavior of the agent is close to
 564 a categorical fifty-fifty evaluation, where all nontrivial uncertainties are weighted the same
 565 without discrimination.

566 The belief component $B(E_a)$ is estimated to be $w^{-1}(W(E_a)) = w^{-1}(0.52) = 0.62$. This
 567 value implies that B must violate additivity. Under additivity, we would have $B(C_a) = 1 -$
 568 $B(E_a) = 0.38$ and then, by symmetry, $B(G_a) = B(S_a) = B(Y_a) = 0.38$, so that $B(G_a) + B(S_a) +$
 569 $B(Y_a) = 3 \times 0.38 = 1.14$. This value should, however, equal $B\{G_a, S_a, Y_a\} = B(E_a)$ under
 570 additivity which is 0.62, leading to a contradiction. Hence, additivity must be violated.

571 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
 572 the result of deviations from risk neutrality that distorted the measurement of $B(E_a)$, and 0.09
 573 is the result of nonadditivity (ambiguity) of belief B . \square

574

575 Theorem 5.2 is valid for virtually all static models of decision under uncertainty and
 576 ambiguity known in the literature today, because Eqs. 5.4 and 5.5 capture virtually all these
 577 models (see Appendix B). Some qualitative observations are as follows. If U is linear, then r
 578 $= w(B(E))$ follows for all $w(B(E)) > 0.5$, providing a very tractable manner of measuring the
 579 nonadditive decision-theory measure $W = w \circ B$.

580 6. Measuring Beliefs through Risk Corrections

581 The next two sections, constituting the second part of the analysis of this paper, analyze
 582 proper scoring rules using the revealed-preference technique. That is, we do not derive

583 empirical predictions from theoretical models, but we reverse the implication. We assume
 584 that empirical observations are given and derive theoretical models from these. In particular,
 585 we will derive beliefs $B(E)$ from reported probabilities r_E . Before turning to this technique,
 586 we discuss alternative measurements of beliefs B considered in the literature.

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

$$589 \quad B(E) = w^{-1}(W(E)). \quad (6.1)$$

590 In general, such revelations of w and W are laborious. The observed choices depend not only
 591 on w and W but also on the utility function U , so that complex multi-parameter estimations
 592 must be carried out (Tversky & Kahneman 1992, p. 311) or elaborate nonparametric
 593 measurements (Abdellaoui, Vossman, & Weber 2005).

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

$$596 \quad (p:x, y) \sim (E:x, y) \quad (6.2)$$

597 for some preset $x > y$, say $x = 100$ and $y = 0$. Then $w(B(E))(U(x)-U(y)) = w(p)(U(x)-U(y))$,
 598 and $B(E) = p$ follows. Wakker (2004) discussed the interpretation of Eqs. 6.1 and 6.2 as
 599 belief. Canonical probabilities were commonly used in early decision analysis (Raiffa 1968,
 600 §5.3; Yates 1990 pp. 25-27) under the assumption of expected utility. A recent experimental
 601 measurement is in Holt (2006, Ch. 30), who also assumed expected utility. Abdellaoui,
 602 Vossman, & Weber (2005) measured and analyzed them in terms of prospect theory, as does
 603 our paper. A practical difficulty is that the measurement of canonical probabilities requires
 604 the measurement of indifferences, and these are not easily inferred from choice. For
 605 example, Holt (2006) used the Becker-deGroot-Marschak mechanism, and Abdellaoui,
 606 Vossman, & Weber (2005) a bisection method. Huck & Weizsäcker (2002) compared the
 607 QSR to the measurement of canonical probabilities and found that the former is more
 608 accurate.

609 A third way to correct reported probabilities is through calibration, where many reported
 610 probabilities are collected over time and then are related to observed relative frequencies.
 611 Calibration has been studied in theoretical game theory (Sandroni, Smorodinsky, & Vohra
 612 2003), and has been applied to weather forecasters (Murphy & Winkler 1974). It needs
 613 extensive data, which is especially difficult to obtain for rare events such as earthquakes, and
 614 further assumptions such as stability over time. Clemen & Lichtendahl (2005) discussed

615 these drawbacks and proposed correction techniques for probability estimates in the spirit of
 616 our paper, but still based these on traditional calibration techniques. Our correction
 617 (“calibration”) technique is considerably more efficient than traditional ones. It shares with
 618 Prelec’s (2004) method the advantage that we need not wait until the truth or untruth of
 619 uncertain events has been revealed for implementing it.

620 We now use the revealed-preference technique to introduce risk corrections. These
 621 combine the advantages of measuring $B(E) = w^{-1}(W(E))$, of measuring canonical
 622 probabilities, and of calibrating reported probabilities relative to objective probabilities, while
 623 avoiding the problems described above, by benefiting from the efficiency of proper scoring
 624 rules. The QSR does entail a restriction of the observations regarding $B(E)$ to cases of E
 625 being more likely than E^c (Observation 4.3.1). The first results do express beliefs B (or p) in
 626 terms of observed values r , but are not complete revealed-preference results because the
 627 right-hand sides of the equations still contain utilities, which are theoretical quantities that are
 628 not directly observable. A “coincidental” agreement of two right-hand sides along the way
 629 will then lead to the main result of this paper: A complete revealed-preference result,
 630 deriving beliefs B entirely from observable choice.

631 We first consider expected utility of §4.2. The following result follows from Theorem
 632 4.2.1 through algebraic manipulations or, for $r \neq 0.5$, as a corollary of Corollary 6.2 hereafter.

633

634 COROLLARY 6.1. Under expected utility with p the (objective or subjective) probability of
 635 event E , and $r = R(p)$ the optimal choice, we have

$$636 \quad p = \frac{r}{r + (1-r) \frac{U'(1-(1-r)^2)}{U'(1-r^2)}} \quad (6.3)$$

637 □

638 We next consider nonexpected utility for known probabilities as in §4.3. An explicit
 639 expression of p in terms of $(U$ and) r , i.e. of $R^{-1}(p)$, follows next for $r > 0.5$, assuming that we
 640 can invert the probability weighting function w . The result follows from Theorem 4.3.3.

641

642 COROLLARY 6.2. Under nonexpected utility with given probabilities (Eq. 4.3.1), the optimal
 643 choice $r = R(p)$ satisfies:

644 If $r > 0.5$, then $p = R^{-1}(r) = w^{-1} \left(\frac{r}{r + (1-r) \frac{U'(1-(1-r)^2)}{U'(1-r^2)}} \right)$. (6.4)

645 □

646 In general, it may not be possible to derive both w and U from $R(p)$ without further
 647 assumptions, i.e. U and w may be nonidentifiable for proper scoring rules. Under regular
 648 assumptions about U and w , however, they have some different implications. The main
 649 difference is that, if we assume that U is differentiable (as done throughout this paper) and
 650 concave, then a flat part of $R(p)$ around 0.5 must be caused by w (Observation 4.2.3).

651 We, finally, turn to nonexpected utility if no probabilities are known, as in §5. Theorem
 652 5.2 implies the following results. It illustrates once more how deviations from expected
 653 utility (w) and nonlinear utility (the marginal-utility ratio) distort the classical proper-scoring-
 654 rule assumption of $B(E) = r$.

655

656 COROLLARY 6.3. Under nonexpected utility with unknown probabilities (Eq. 5.4), the
 657 optimal choice $r = r_E$ satisfies:

658 If $r > 0.5$, then $B(E) = w^{-1} \left(\frac{r}{r + (1-r) \frac{U'(1-(1-r)^2)}{U'(1-r^2)}} \right)$. (6.5)

659 □

660

661 As preparation for a complete revealed-preference result, note that the right-hand sides
 662 of Eqs. 6.4 and 6.5 are identical. Hence, if we find a p in Eq. 6.4 with the same r value as E ,
 663 then we can, by Eq. 6.4, immediately substitute p for the right-hand side of Eq. 6.5, getting
 664 $B(E) = p$ without need to know the ingredients w and U of Eq. 6.5. This observation (to be
 665 combined with Eq. 2.2 for $r < 0.5$) implies the following corollary, which is the main result of
 666 this paper and which is displayed for its empirical importance.

667

668 COROLLARY 6.4. Under nonexpected utility with unknown probabilities (Eq. 5.4), assume for
 669 the optimal choice $r = r_E$ that $r > 0.5$. Then

$$670 \quad B(E) = R^{-1}(r). \quad (6.6)$$

671 \square

672

673 This corollary is useful for empirical purposes. It is the only implication of our
 674 theoretical analysis that is needed for applications. It shows how proper scoring rules can
 675 allow for deviations from expected value and expected utility, and is key in filtering out risk
 676 attitudes. We first infer the (for the participant) optimal $R(p)$ for a set of exogenously given
 677 probabilities p that is so dense (all values $p = j/20$ for $j \geq 10$ in our experiment) that we obtain
 678 a sufficiently accurate estimation of R and R^{-1} . Then, for all uncertain events E more likely
 679 than their complement, we immediately derive $B(E)$ from the observed r_E through Eq. 6.6.

680 Summarizing:

681

682 If for event E the participant reports probability $r_E = r$
 683 and for objective probability p the participant also reports probability $R(p) = r$
 684 then $B(E) = p$.

685

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

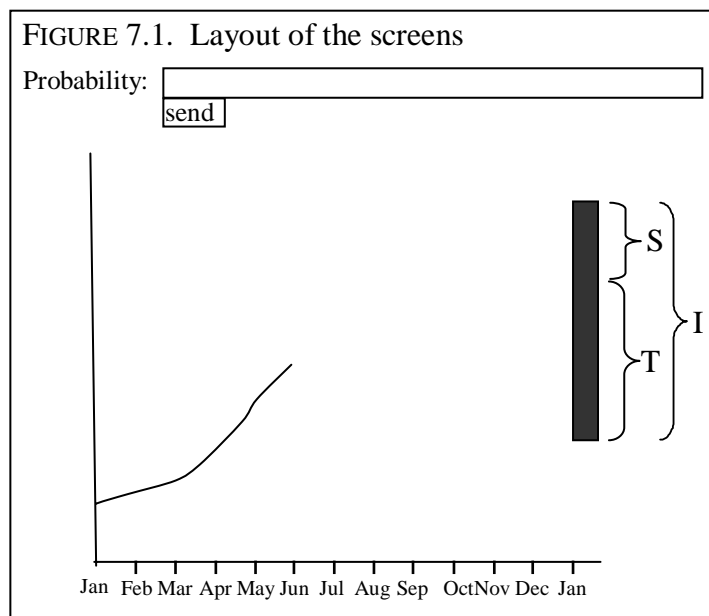
689 We call the function R^{-1} the *risk-correction* (for proper scoring rules), and $R^{-1}(r_E)$ the
 690 *risk-corrected probability*. This value is the canonical probability, obtained without having
 691 measured indifference curves such as through the Becker-DeGroot-Marschak mechanism, without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

697 The curves in Figure 3.1 can be reinterpreted as inverses of risk corrections. The
 698 examples illustrated there were based on risk averse decision attitudes, leading to
 699 conservative estimations moved in the direction of 0.5. Risk seeking will lead to the opposite
 700 effect, and will generate overly extreme reported probabilities, suggesting overconfidence.
 701 Obviously, if factors in the probability elicitation of the calibration part induce
 702 overconfidence and risk seeking, then our risk correction will detect those biases and correct

703 for them. If, after the risk correction, overconfidence is (still) present, then it cannot be due
 704 to risk seeking. We can then conclude with more confidence that overconfidence is a
 705 genuine property of belief, irrespective of risk seeking.

706 **7. An Illustration of Our Measurement of Belief**

707 This section describes risk corrections for a participant in the experiment so as to
 708 illustrate how our method can be applied empirically. We will see that Corollary 6.4 is the
 709 only result of the theoretical analysis needed to apply our method. Results and curves for $r <$
 710 0.5 are derived from $r > 0.5$ using Eq. 2.2; we will not mention this point explicitly in what
 711 follows.



724 The left side of Figure 7.1 displays the performance of stock 12 in our experiment from
 725 January 1 until June 1 1991 as given to the participants. Stock 12 concerned the stock
 726 Begemann Kon. Groep (General Industries). Further details (such as the absence of a unit on
 727 the y-axis) will be explained in §8. The right side of the figure displays two disjoint intervals
 728 S and T , and their union $I = S \cup T$. For each of the intervals S, T , and I , participants reported
 729 the probability of the stock ending up in that interval on January 1 1992 (with some other
 730 questions in between these three questions). For participant 14, the results are as follows.

731 $r_S = 0.35; r_T = 0.55; r_I = 0.65.$ (7.1)

732 Under additivity of reported probability, $r_S + r_T - r_I$ (the *additivity bias*, defined in general in
733 Eq. 8.5), should be 0, but here it is not and additivity is violated.

734 The additivity bias is $0.35 + 0.55 - 0.65 = 0.25$. (7.2)

735 Table 7.1 and Figure 7.2 (in inverted form) display the reported probabilities $R(p)$ that
736 we measured from this participant, with the curves explained later. We use progressive
737 averages (midpoints between data points) so as to reduce noise.

738

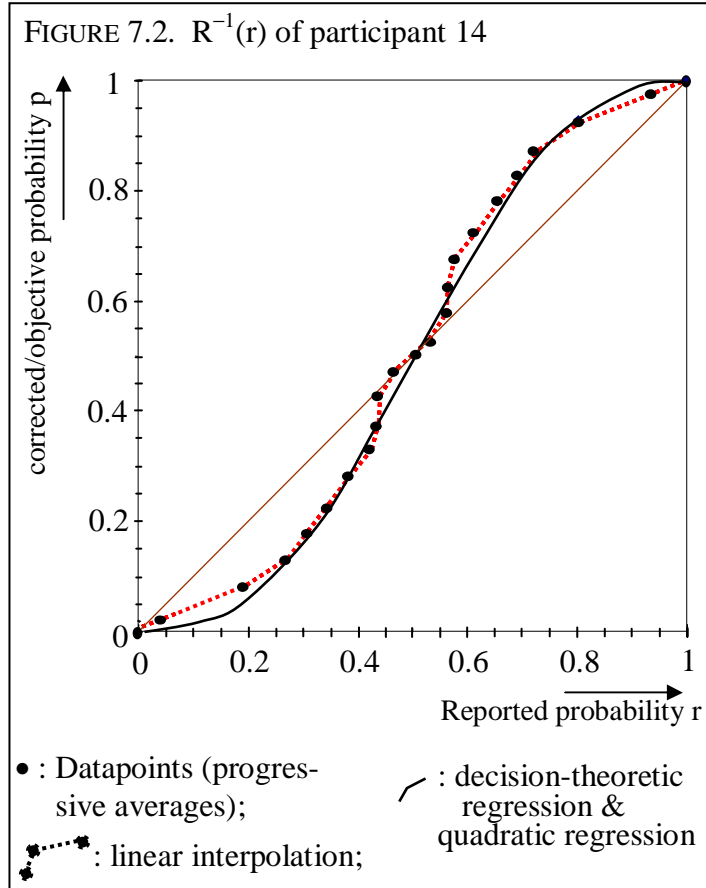
TABLE 7.1. Progressive average reported probabilities $R(p)$ of participant 14

p	.025	.075	.125	.175	.225	.275	.325	.375	.425	.475	.525	.575	.625	.675	.725	.775	.825	.875	.925	.975
$R(p)$.067	.192	.267	.305	.345	.382	.422	.435	.437	.470	.530	.563	.565	.578	.618	.655	.695	.733	.808	.933

739

740

741



742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

For simplicity of presentation, we analyze the data here using linear interpolation. Then $R(0.23) = 0.35$.⁶ Using this value for $R(0.23)$, using the values $R(0.56) = 0.55$ and $R(0.77) = 0.65$, and, finally, using Eq. 6.6, we obtain the following risk-corrected beliefs.

$$\begin{aligned}
 B(S) &= R^{-1}(0.35) = 0.23; \quad B(T) = R^{-1}(0.55) = 0.56; \quad B(I) = R^{-1}(0.65) = 0.77; \\
 \text{the additivity bias is } &0.23 + 0.56 - 0.77 = 0.02.
 \end{aligned}
 \tag{7.3}$$

⁶ We have $0.23 = 0.865 \times 0.225 + 0.135 \times 0.275$, $R(0.225) = 0.345$, and $R(0.275) = 0.382$, so that $R(0.23) = R(0.865 \times 0.225 + 0.135 \times 0.275) = 0.865 \times R(0.225) + 0.135 \times R(0.275) = 0.865 \times 0.345 + 0.135 \times 0.382 = 0.35$.

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

769 It is statistically preferable to fit data with smoother curves than resulting from linear
 770 interpolation. We derived “decision-theoretic” parametric curves for $R(p)$ from Corollary
 771 6.2, with further assumptions explained at the end of §9.1.⁷ The resulting curve for
 772 participant 14 is given in the figure. The equality $B = R^{-1}(r)$ and this curve lead to

$$773 \quad B(S) = R^{-1}(0.35) = 0.24; B(T) = R^{-1}(0.55) = 0.59; B(I) = R^{-1}(0.65) = 0.76; \text{ the additivity} \\
 774 \quad \text{bias is } 0.24 + 0.59 - 0.76 = 0.07, \quad (7.4)$$

775 again reducing the uncorrected additivity bias. The quadratic curve, explained in §11, is
 776 indistinguishable from the decision theoretic curve.

777

778 **8. An Experimental Application of Risk Corrections: Method**

779 The following four sections present the third part of this paper, being an experimental
 780 implementation of our new measurement method.

781

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

785

786 *Procedure.* Participants were seated in front of personal computers in 6 groups of
 787 approximately 16 participants each. They first received an explanation of the QSR, given in
 788 Appendix C. Then, for each uncertain event, participants could first report a probability (in

⁷ The decision-theoretic curve in the figure is the function $p = B(E) = \frac{r}{r + (1-r) \frac{-0.26(1-(1-r)^2)^{-1.26}}{-0.26(1-r^2)^{-1.26}}}$, in

agreement with Corollaries 5.4 and 5.2, where we estimated $w(p) = p$ and found $\rho = -0.26$ as optimal value for $U(x)$ in Eq. 8.1.

789 percentages) by typing in an integer from 0 to 100. Subsequently, the confirmation screen
 790 displayed a list box with probabilities and the corresponding score when the event was (not)
 791 true, illustrated in Figure 8.1.

792

793

794

795

796

797

798

799

800

801

802

803

804

FIGURE 8.1.

Probability	Your score if statement is true	Your score if statement is not true
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

send

805 All figures (including Figure 7.1) are reproduced here in black and white; in the experiment
 806 we used colors to further clarify the figures. The entered probability and the corresponding
 807 score were preselected in this list box. The participant could confirm the decision or change
 808 to another probability by using the up or down arrow or by scrolling to another probability
 809 using the mouse. The event itself was also visible on the confirmation screen. Thus, the
 810 reported probability r finally resulted for the uncertain event.

811

812 *Stimuli*

813 The participants provided 100 reported probabilities r for events with unknown probabilities
 814 in the *stock-price part* of the experiment. For these events, we fixed June 1, 1991, as the
 815 “evaluation date.” The uncertain events always concerned the question whether or not the
 816 price of a stock would lie in a target-interval seven months after the evaluation date. For
 817 each stock, the participants received a graph depicting the price of the stock on 0, 1, 2, 3, 4,
 818 and 5 months prior to the evaluation date, as well as an upper and lower bound to the price of
 819 the stock on the evaluation date. Figure 7.1, without the braces and letters, gives an example
 820 of the layout. We used 32 different stocks, all real-world stock market data from the 1991
 821 Amsterdam Stock Exchange. After 4 practice questions, the graph of each stock-price was
 822 displayed once in the questions 5-36, once in the questions 37-68, and once in the questions

823 69-100. We, thus, obtained three probabilistic judgments of the performance of each stock,
824 once for a large target-interval and twice for small target-intervals that partitioned the large
825 target-interval (see Figure 7.1). We partially randomized the order of presentation of the
826 elicitations. Each stock was presented at the same place in the first, second, and third 32-
827 tuple of elicitations, so as to ensure that questions pertaining to the same stock were always
828 far apart. The order of presentation of the two small and one large interval for each stock
829 were not randomized stochastically, but were varied systematically, so that all orders of big
830 and small intervals occurred equally often. We also maximized the variation of whether
831 small intervals were both very small, both moderately small, or one very small and one
832 moderately small.

833

834 In the *calibration part* of the experiment, participants made essentially the same decisions as
835 in the stock-price part, but now for 20 events with objective probabilities. Thus, participants
836 simply made choices between risky prospects with objective probabilities. We used two 10-
837 sided dice to determine the outcome of the different prospects and obtained measurements of
838 the reported probabilities corresponding to the objective probabilities 0.05, 0.10, 0.15, ...,
839 0.85, 0.90, and 0.95 (we measured the objective probability 0.95 twice). The event with
840 probability 0.25 was, for instance, described as “The outcome of the roll with two 10-sided
841 dice is in the range 01–25.” The decision screen was very similar to Figure 8.1, except for
842 the fact that we wrote “row-percentage” instead of “probability” and “your score if the roll of
843 the die is 01-25” instead of “your score if statement is true;” etc.

844

845 *Motivating participants.* Depending on whether the uncertain event obtained or not and on
846 the reported probability for the uncertain event, a number of points was determined for each
847 question through the QSR (Eq. 2.1), using 10000 points as unit of payment so as to have
848 integer scores with four digits of precision. Thus, the maximum score for one question was
849 10000, the minimum score was 0, and the certain score resulting from reported probability
850 0.5 was 7500 points.

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

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

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

861

862 *Analysis.* For the calibration part we only need to analyze probabilities of 0.5 or higher, by
 863 Observation 4.3.2. Indeed, by Eq. 3.2, every observation for $p < 0.5$ amounts to an
 864 observation for $p' = 1-p > 0.5$. It implies that we have two observations for all $p > 0.5$ (and
 865 three for $p = 0.95$).

866 We first analyze the data at the group level, assuming homogeneous participants. We
 867 start from the general model of Eq. 4.3.1. Notice that this equation can be estimated using a
 868 non-parametric procedure. If the agent is willing to go through a large series of correction
 869 questions, it is possible to measure the corresponding reported probability of each objective
 870 probability repeatedly. In this way an accurate estimate of the whole correction curve can be
 871 obtained without making assumptions about the utility function or the weighting function.
 872 This procedure seems the appropriate one if the goal is to correct an expert, e.g., correct the
 873 reports provided by a weatherman. In applications of experimental economics where
 874 subjects participate for a limited amount of time, the researcher will only be able to collect a
 875 limited number of observations of the correction curve. Then it is more appropriate to follow
 876 a parametric approach to elicit the curve that fits the observations best. In this paper, we used
 877 parametric fittings. For U we used the *power utility with parameter* ρ , also known as the
 878 family of constant relative risk aversion (CRRA)⁸, and the most popular parametric family
 879 for fitting utility, which is defined as follows:

$$\begin{aligned}
 880 \quad & \text{For } \rho > 0: U(x) = x^\rho; \\
 881 \quad & \text{for } \rho = 0: U(x) = \ln(x); \\
 882 \quad & \text{for } \rho < 0: U(x) = -x^\rho.
 \end{aligned}
 \tag{8.1}$$

883 It is well-known that the unit of payment is immaterial for this family. The most general
 884 family that we consider for $w(p)$ is Prelec's (1998) 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.

885 $w(p) = \left(\exp(-\beta(-\ln(p))^\alpha) \right),$ (8.2)

886 chosen for its analytic tractability and good empirical performance. We will mostly use the
 887 one-parameter subfamily with $\beta=1$, as in Eq. 4.3.4, for reasons explained later. Substituting
 888 the above functions yields

889
$$B(E) = \exp\left(-\left(\frac{-\ln\left(\frac{r(2r-r^2)^{1-p}}{(1-r)(1-r^2)^{1-p} + r(2r-r^2)^{1-p}}\right)}{\beta}\right)^{1/\alpha}\right).$$

890 for Eq. 6.5.

891 The model we estimate is as follows.

892 $R_{s,t,k}(j/20) = h(j/20, \alpha_t, \rho_t) + \varepsilon_{s,t,k}(j/20, \sigma_t^2).$ (8.3)

893 Here $R_{s,t,k}(j/20)$ is the reported probability of participant s for known probability $p=j/20$ ($10 \leq$
 894 $j \leq 19$) in treatment t ($t = \text{ALL}$ or $t = \text{ONE}$) for the k^{th} measurement for this probability, with
 895 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
 896 is the remaining probability-weighting parameter (Eq. 8.2), and ρ_t is the power of utility (Eq.
 897 8.1). The function h is the inverse of Eq. 6.4. Although we have no analytic expression for
 898 this inverse, we could calculate it numerically in the analyses. The error terms $\varepsilon_{s,t,k}(j/20)$ are
 899 drawn from a truncated normal distribution with mean 0 and treatment dependent variance
 900 σ_t^2 . The distribution of the error terms is truncated because reported probabilities below 0
 901 and above 1 are excluded by design. Error terms are identically and independently
 902 distributed across participants and choices. We employed maximum likelihood to estimate
 903 the parameters of Eq. 8.3. We also carried out an analysis at the individual level of the
 904 calibration part, with $\alpha_{s,t}$ and $\rho_{s,t}$ instead of α_t and ρ_t , i.e. with these parameters depending on
 905 the participant.

906 In the stock-price part, violations of additivity were tested. With I the large interval of a
 907 stock, being the union $S \cup T$ of the two small intervals S and T , additivity of the uncorrected
 908 reported probabilities implies

909 $r_S + r_T = r_I.$ (8.4)

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

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

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

$$917 \quad B(S) + B(T) = B(I). \quad (8.5)$$

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

920 We next discuss tests of the additivity bias. For each individual stock, and also for the
 921 average over all stocks, we tested for both treatments $t=ONE$ and $t=ALL$: (a) whether the
 922 additivity bias was zero or not, both with and without risk correction; (b) whether the average
 923 additivity bias, as relevant for aggregated group behavior and expert opinions, was enlarged
 924 or reduced by correction; (c) whether the absolute value of the additivity bias, as relevant for
 925 additivity at the individual level, was enlarged or reduced by correction. We report only the
 926 tests for averages over all stocks.

927

928 **9. Results of the Calibration Part**

929 Risk-corrections and, in general, QSR measurements, do not make sense for participants who
 930 are hardly responsive to probabilities, so that $R(p)$ is almost flat on its entire domain. Hence
 931 we kept only those participants for whom the correlation between reported probability and
 932 objective probability exceeded 0.2. We thus dropped 4 participants. The following analyses
 933 are based on the remaining 89 participants.

934

935 *9.1. Group Averages*

936

937 We did several tests using Eq. 8.2 with β as a free (treatment-dependent or -independent)
 938 variable, but β 's estimates added little extra explanatory power to the other parameters and
 939 usually were close to 1. Hence, we chose to focus on a more parsimonious model in which
 940 the restriction $\beta_{ONE} = \beta_{ALL} = 1$ is employed. Table 9.1 lists the estimates for the model of Eq.
 941 8.3 for $\beta=1$ (Eq. 4.3.4 instead of Eq. 8.2) together with the estimates of some models with

942 additional restrictions. We first give results for group averages, assuming homogeneous
943 participants.

944

945 *Overall need for risk-correction.* The 1st row of Table 9.1 shows the results for the most
946 general model. The 2nd row presents the results without any correction. The likelihood
947 reduces significantly (Likelihood Ratio test, $p = 0.01$) and substantially, so that risk-
948 correction is called for. Risk-correction is also called for in both treatments in isolation, as
949 the 3rd and 4th rows show, which significantly improve the likelihood relative to the 2nd row
950 (Likelihood Ratio test; $p = 0.01$ for $t=ALL$, comparing 3rd to 2nd row; $p = 0.01$ for $t=ONE$,
951 comparing 4th to 2nd row).

952

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

961 TABLE 9.1. Estimation results at the aggregate level

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

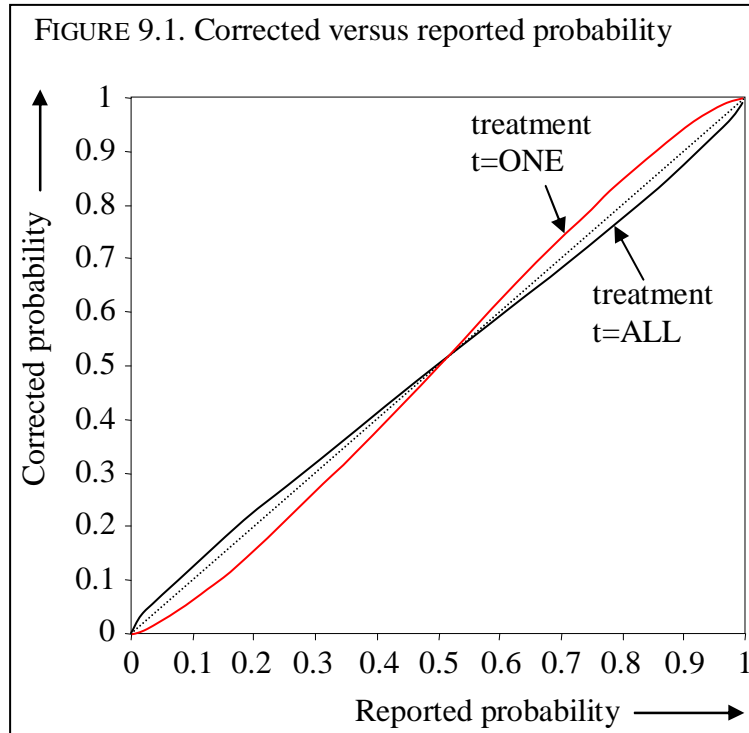
962 Standard errors in parentheses, * denotes significance at the 1% level.

963

964 The difference between the two treatments is apparently caused by curvature of utility,
 965 captured by ρ_{ONE} and ρ_{ALL} . We obtain $\rho_{ONE} < \rho_{ALL}$: when only one decision is paid out then
 966 participants exhibit more concave curvature of utility than when all decisions are paid out.
 967 Given same probability weighting, it implies more risk aversion for t=ONE than for t=ALL
 968 (and R closer to 0.5). The finding is supported by comparing the 6th row of Table 9.1, with
 969 the restriction $\rho_{ONE} = \rho_{ALL}$, to the 1st row. This restriction significantly reduces the
 970 likelihood of observing the data (Likelihood Ratio test, $p = 0.01$).

971

972 *Comparing utility and probability weighting.* Correcting only for utility curvature (7th row,
 973 $\alpha_{\text{ONE}} = \alpha_{\text{ALL}} = 1$) has a somewhat better likelihood than correcting only for probability
 974 weighting (8th row, $\rho_{\text{ONE}} = \rho_{\text{ALL}} = 1$).



990

991 *Discussion of comparison of utility curvature and probability weighting for group-averages.*

992 In deterministic choice, α could be determined through the flat part of R around 0.5, after
 993 which ρ could serve to improve the fit elsewhere. Statistically, however, α and ρ have much
 994 overlap, with risk aversion enhanced and $R(p)$ moved towards 0.5 by increasing α and
 995 decreasing ρ , and one does not add much explanatory power to the other. It is, therefore,
 996 better to use only one of these parameters. Another reason to use only one parameter
 997 concerns the individual analysis reported in the following subsection. Because we only have
 998 20 choices per participant it is important to economize on the number of free parameters
 999 there.

1000 We found that ρ has a slightly better explanatory power than α . For this reason, and for
 1001 reasons of convenience (see discussion section), we will only use the parameter ρ , and
 1002 assume $\alpha = 1$ henceforth. Figure 9.1 displays the resulting average risk-correction for the
 1003 two treatments separately.

1004

1005 *Comparing the two treatments when there is no probability weighting.* The average effect of
 1006 correction for utility curvature is not strong, especially for t=ALL. Yet this correction has a
 1007 significant effect, as can be seen from comparing the 7th row (general ρ) in Table 9.1 to its 9th
 1008 row ($\rho_{ALL} = 1$) (Likelihood Ratio test, $p = 0.01$).

1009

1010

9.2. Individual Analyses

1011

1012 *Need for risk-correction at the individual level.* There is considerable heterogeneity in each
 1013 treatment. Whereas the corrections required were significant but small at the level of group
 1014 averages, they are big at the individual level. This appears from Figure 9.2, which displays
 1015 the cumulative distribution of the (per-subject) estimated ρ -coefficients for each treatment,
 1016 assuming $\alpha = \beta = 1$. There are wide deviations from the value $\rho=1$ (i.e., no correction) on
 1017 both sides. As seen from the group-average analysis, there are more deviations at the risk-
 1018 averse side of $\rho < 1$.

1019

1020

1021

1022

1023

1024

1025

1026

1027

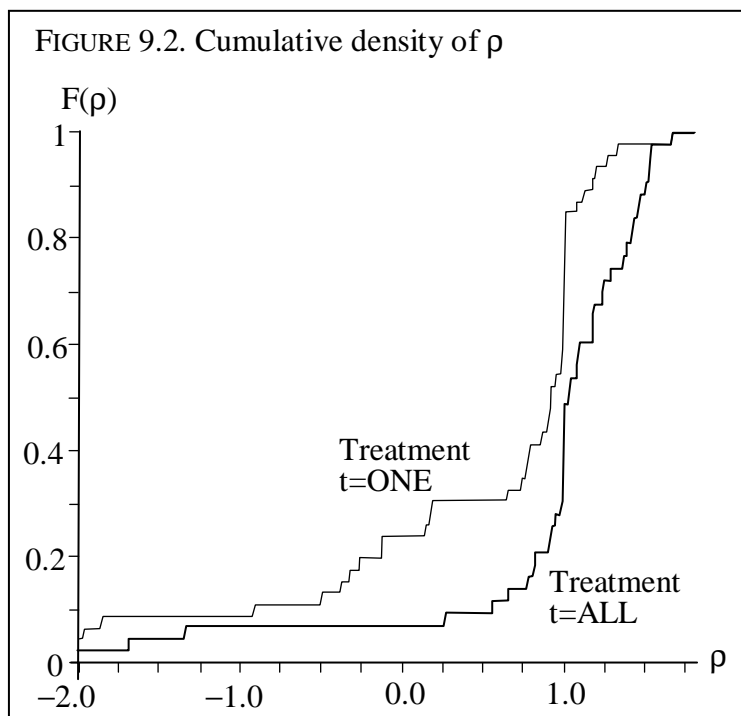
1028

1029

1030

1031

1032



1033

1034 *Comparing the two treatments.* The ρ -coefficient distribution of treatment t=ONE dominates
 1035 the ρ -coefficient distribution of treatment t=ALL. A two-sided Mann-Whitney test rejects
 1036 the null-hypothesis that the ranks of ρ -coefficients are equal across the treatments in favor of
 1037 the hypothesis that the ρ -coefficients for t=ONE are lower than for t=ALL ($p=0.001$). It

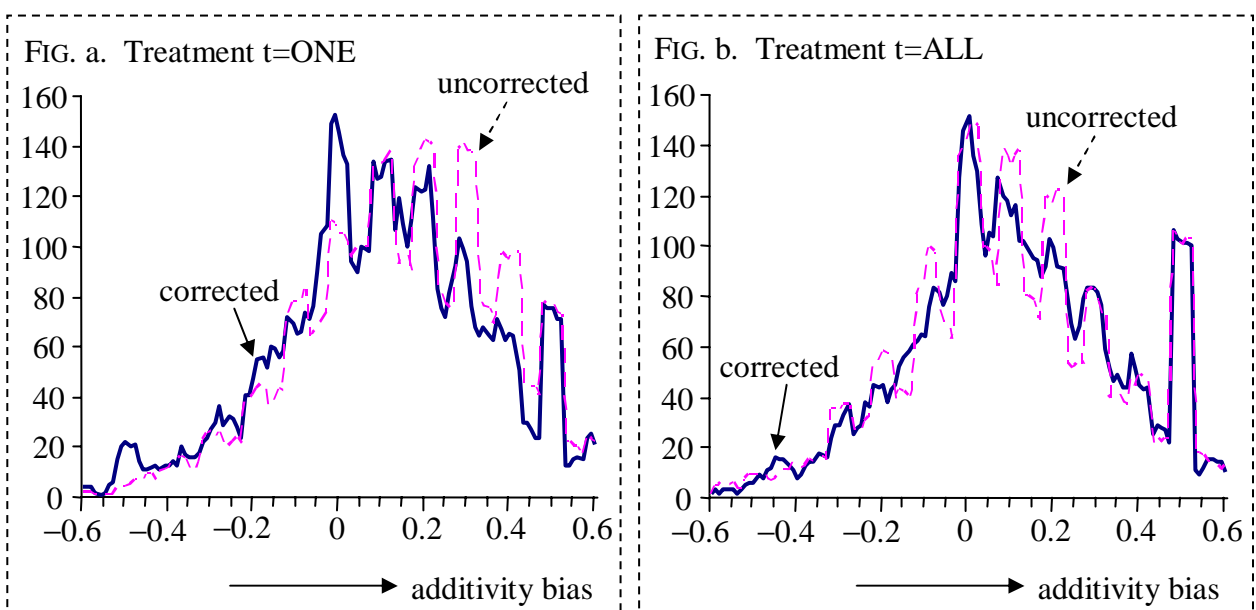
1038 confirms that for group averages there is more risk aversion, moving R in the direction of 0.5,
1039 for t=ONE than for t=ALL. The figure also shows that in an absolute sense there is more
1040 deviation from $\rho=1$ for t=ONE than for t=ALL, implying that there are more deviations from
1041 expected value and more risk corrections for t=ONE than for t=ALL.

1042

1043 Unlike the median ρ -coefficients that are fairly close to each other for the two treatments
1044 (0.92 for t=ONE versus 1.04 for t=ALL), the mean ρ -coefficients are substantially different
1045 (0.24 for t=ONE versus 0.91 for t=ALL), which is caused by skewedness to the left for
1046 t=ONE. That is, there is a relatively high number of strongly risk-averse participants for
1047 t=ONE. Analyses of the individual ρ parameters (two-sided Wilcoxon signed rank sum tests)
1048 confirm findings of group-average analyses in the sense that the ρ -coefficients are
1049 significantly smaller than 1 for t=ONE ($z = -3.50$, $p = 0.0005$), but not for t=ALL ($z = 1.42$,
1050 $p = 0.16$).

1051 **10. Results for the Stock-Price Part: Risk-Correction and** 1052 **Additivity**

FIGURE 10.1. Empirical density of additivity bias for the two treatments



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 additivity biases in the interval, aggregated over 32 stocks and 89 individuals, for both treatments. With risk-correction, there were 65 additivity biases between 0.375 and 0.425 in the treatment t=ONE, and without risk-correction there were 95 such; etc.

All comparisons hereafter are based on two-sided Wilcoxon signed rank sum tests. Figure 10.1 displays data, aggregated over both stocks and individuals, of the additivity biases for t=ONE and for t=ALL. The figures show that the additivity bias is more often positive than negative, in agreement with common findings in the literature (Tversky & Koehler 1994; Bateman et al. 1997). Indeed, for virtually all stocks the additivity bias is significantly positive for both treatments, showing in particular that additivity does not hold. This also holds when 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 corrections reduce the violations of additivity.

Let us first consider t=ONE. Here the risk corrections reduce the average additivity bias significantly for 27 of the 32 stocks, and enlarge it for none. We only report the statistics for the average additivity bias over all stocks per individual, which has overall averages 0.163 (uncorrected) and 0.120 (corrected), with the latter significantly smaller ($z = 3.21$, $p = 0.001$). For assessing the degree of irrationality (additivity-violation) at the individual level, the absolute values of the additivity bias are interesting. For t=ONE, Figure 10.1 suggests that these are smaller after correction, because on average the corrected curve is closer to 0 on the

1087 x-axis. These absolute values were significantly reduced for 9 stocks and enlarged for none.
1088 Again, we only report the statistics for the average absolute value of the additivity bias over
1089 all stocks per individual, which has overall averages 0.239 (uncorrected) and 0.228
1090 (corrected), with the latter significantly smaller ($z = 2.26$, $p = 0.02$).

1091 For $t=ALL$, risk corrections did not significantly alter the average additivity bias. More
1092 precisely, it gave a significant increase for 3 stocks and a significant decrease for 1 stock,
1093 which, for 32 stocks, suggests no systematic effect. The latter was confirmed when we took
1094 the average additivity bias over all stocks for each individual, with no significant differences
1095 generated by correction (average 0.128 uncorrected and average 0.136 corrected; $z = -1.64$, p
1096 $= 0.1$). Similar results hold for absolute values of additivity biases, which gave a significant
1097 increase for 1 stock and a significant decrease for no stock. Taking the average additivity
1098 bias over all stocks for each individual (average 0.237 uncorrected and average 0.239
1099 corrected; $z = -0.36$, $p = 0.7$) also gave no significant difference.

1100 Classifications of individuals according to whether they exhibited more positive or more
1101 negative additivity biases, and to whether risk corrections improved or worsened the
1102 additivity bias more often, confirmed the patterns obtained above through stockwise
1103 analyses, and will not be reported.

1104 Risk correction reduces the additivity bias for treatment $t=ONE$ to a level similar to that
1105 observed for $t=ALL$ (averages 0.120 and 0.136). The overall pattern is that beliefs for
1106 $t=ONE$ after correction, and for $t=ALL$ both before and after correction, exhibit a similar
1107 degree of violation of additivity, which is clearly different from zero. The additivity bias is
1108 not completely caused by nonlinear risk attitudes when participants report probabilities, but
1109 has a genuine basis in beliefs.

1110

1111 **11. Discussion of Experiment**

1112 *Methods.* We chose the evaluation date (June 1, 1991) sufficiently long ago to ensure that
1113 participants would be unlikely to recognize the stocks or have private information about
1114 them. In addition, no numbers were displayed on the vertical axis, making it extra hard for
1115 participants to recognize specific stocks. We, thus, ensured that participants based their
1116 probability judgments entirely on the prior information about past performance of the stocks
1117 given by us. Given the large number of questions it is unlikely that participants noticed that

1118 the graphs were presented more than once (three times) for each stock. Indeed, in informal
1119 discussions after the experiment no participant showed awareness of this point.

1120 In some studies in the literature, the properness of scoring rules is explained to
1121 participants by stating that it is in their best interest to state their true beliefs, either without
1122 further explanation, or with the claim added that they will thus maximize their “expected”
1123 money. A drawback of this explanation is that expected value maximization is empirically
1124 violated, which is the central topic of this paper (§3), so that the recommendation is
1125 debatable. We, therefore, used an alternative explanation that relates properness for one-off
1126 events to observed frequencies of repeated events (Appendix C).

1127

1128 *Optimal Incentive Scheme.* After some theoretical debates about the random-lottery incentive
1129 system (Holt 1986), as in our treatment t=ONE, the system was tested empirically and found
1130 to be incentive-compatible (Starmer & Sugden 1991). It is today the almost exclusively used
1131 incentive system for measurements of individual preferences (Holt & Laury 2002; Harrison
1132 et al. 2002). Unlike repeated payments it avoids income effects such as Thaler & Johnson's
1133 (1990) house money effect, and the drift towards expected value and linear utility that is
1134 commonly generated by repeated choice.⁹ For the purpose of measuring individual
1135 preference, the treatment t=ONE is, therefore, preferable. When the purpose is, however, to
1136 derive subjective probabilities from proper scoring rules, and no risk-correction is possible,
1137 then a drift towards expected value is actually an advantage, because uncorrected proper
1138 scoring rules assume expected value. This point agrees with our findings, where less risk-
1139 correction was required for the t=ALL treatment. Li (2007) discusses other arguments for
1140 and against repeated rewarding when events are not verifiable and when binary rewards have
1141 to be used.

1142 For some applications group averages of probability estimates are most relevant, such as
1143 when aggregating expert judgments or predicting group behavior. Then our statistical results
1144 regarding “non-absolute” values of reported probabilities are most relevant. For the
1145 assessment of rationality at the individual level, absolute values of the additivity biases are
1146 most relevant.

1147

⁹ It is required that the repeated choices are perceived as sufficiently uncorrelated. Correlation can enhance the perception of and aversion to ambiguity (Halevy & Feltkamp 2005).

1148 *Choice of Parameters.* The lack of extra explanatory power of parameter β in Eq. 8.2 should
1149 come as no surprise because β and α behave similarly on $[0.5,1]$, increasing risk aversion
1150 there. They mainly deviate from one another on $[0,0.5]$, where β continues to enhance risk
1151 aversion but α enhances the inverse-S shape that is mostly found empirically. The domain
1152 $[0,0.5]$ is, however, not relevant to our study (Observation 4.3.2).

1153 We found that the risk correction through the utility curvature parameter ρ fitted the data
1154 somewhat better than the correction through the probability-weighting parameter α . This
1155 finding may be interpreted as some descriptive support for expected utility. Another reason
1156 that we used ρ and not α in our main analysis is that ρ , and utility curvature, are more well-
1157 known in the economic literature than probability weighting, and are more analytically
1158 tractable with R^{-1} defined everywhere. Although ρ indeed reflects the power of utility *if*
1159 *expected utility is assumed*, we caution against unqualified interpretations here, as in any
1160 study of risk aversion. The parameter ρ may also capture risk aversion generated by
1161 probability weighting, and possibly by other factors.

1162
1163 *Pragmatic applications.* More tractable families can be used to fit the reported probabilities
1164 than the decision-theory-based curves that we used. For example, in Figure 7.2 we also used
1165 quadratic regression to find the curve $p = a + br + cr^2$ that best fits the data. The curve is
1166 virtually indistinguishable from the decision-theoretic curve. This observation, together with
1167 Corollary 6.4 demonstrating that we only need the readily observable reported probabilities
1168 and not the actual utility function or probability weighting function to apply our method,
1169 shows that applications of our method are easy. The theoretical analysis of this paper, and
1170 the decision-theory based curve-fitting that we adopted, served to prove that our method is in
1171 agreement with modern decision theories. If this thesis is accepted, and the only goal is to
1172 obtain risk-corrected reported probabilities, then one may choose the pragmatic shortcuts just
1173 described.

1174
1175 *General Discussion.* Under proper scoring rules, beliefs are derived solely from decisions,
1176 and Eq. 2.1 is taken purely as a decision problem, where the only goal of the agent is to
1177 optimize the prospect received. Thus, this paper has analyzed proper scoring rules purely
1178 from the decision-theoretic perspective supported with real incentives, and has corrected only
1179 for biases resulting therefrom. Many studies have investigated direct judgments of belief

1180 without real incentives, and then many other aspects play a role, leading for instance to the
1181 often found overconfidence. Such introspective effects are beyond the scope of this paper.

1182 The experimental data show that for a subset of the subjects a substantial correction of
1183 reported probabilities needs to be made. The fraction of the population that needs substantial
1184 corrections is larger when only one big decision is paid than when repeated small decisions
1185 are paid. Our conclusion is that it is desirable to correct agents' reported probabilities elicited
1186 with scoring rules, especially if only a single large decision is paid. If it is not possible to
1187 obtain individual measurements of the correction curve, then it will be useful to use best-
1188 guess corrections, for instance through averages obtained from individuals as similar as
1189 possible. Thus, at least, the systematic error for the group average to risk attitude has been
1190 corrected for as good as is possible without requiring extra measurements. In this respect the
1191 average curves in our Figure 9.1 are reassuring for existing studies, because these curves
1192 suggest that only small corrections were called for regarding the group averages in our
1193 context.

1194 Allen (1987) proposed to avoid biases of the QSR resulting from nonlinear utility by
1195 paying in terms of the probability of winning a prize instead of in terms of money, and this
1196 procedure was implemented by McKelvey & Page (1990). The procedure, however, only
1197 works if expected utility holds, and there is much evidence against this assumption. Indeed,
1198 Selten, Sadrieh, & Abbink (1999) showed empirically that payment in probability does
1199 enhance the desired risk neutral behavior.

1200 **12. Conclusion**

1201 This paper has applied modern theories of risk and ambiguity to proper scoring rules.
1202 Mutual benefits have resulted for practitioners of proper scoring rules and for the study of
1203 risk and ambiguity. For the former we have shown which distortions affect their common
1204 measurements and how large these distortions are, using theories that are descriptively better
1205 than the expected value hypothesis still common for proper scoring rules today. We have
1206 provided a procedure to correct for the aforementioned distortions, and a theoretical
1207 foundation has been given for interpretations of the resulting measurements as (possibly non-
1208 Bayesian) beliefs and/or ambiguity attitudes. For studies of risk and ambiguity we have
1209 shown how the remarkable efficiency of proper scoring rules can be used to measure and

1210 analyze subjective beliefs and ambiguity attitudes in ways more tractable than is possible
1211 through the binary preferences traditionally used.

1212 We have demonstrated the feasibility and tractability of our method in an experiment,
1213 where we used it to investigate some properties of beliefs and quadratic proper scoring rules.
1214 We found, for instance, that our correction method reduces the violations of additivity in
1215 subjective beliefs but does not eliminate them. It confirms that beliefs are genuinely non-
1216 Bayesian and that ambiguity attitudes play a central role.
1217

1218 **Appendix A. Proofs and Technical Remarks**

1219 In Eqs. 4.3.1 and 4.3.2, probability p has a different decision weight when it yields the
1220 best outcome of the prospect ($r > 0.5$) than when it yields the worst ($r < 0.5$). Similarly, in
1221 Eqs. 5.4 and 5.5, E has a different decision weight when it yields the highest outcome ($r >$
1222 0.5) than when it yields the lowest outcome ($r < 0.5$). Such a dependency of decision weights
1223 on the ranking position of the outcome is called *rank-dependence* in the literature.

1224 Under rank-dependence, the sum of the decision weights in the evaluation of a prospect
1225 are 1 even though $w(B(E))$ is not additive in E . This property is necessary for the functional
1226 that evaluates prospects to satisfy natural conditions such as stochastic dominance, which
1227 explains why theoretically sound nonexpected utility models could only be developed after
1228 the discovery of rank dependence, a discovery that was made independently by Quiggin
1229 (1982) for the special case of risk and by Schmeidler (1989, first version 1982) for the
1230 general context of uncertainty.

1231 For qsr-prospects in Eq. 2.1, every choice $r < 0$ is inferior to $r = 0$, and $r > 1$ is inferior to
1232 $r = 1$. The optimization problem does not change if we allow all real r , instead of $0 \leq r \leq 1$.
1233 Hence, solutions $r = 0$ or $r = 1$ hereafter can be treated as interior solutions, and they satisfy
1234 the first-order optimality conditions.

1235

1236 PROOF OF OBSERVATION 4.2.3. If $r = 0.5$ then the marginal utility ratio in Eq. 4.2.2 is 1, and
1237 $p = 0.5$ follows. For the reversed implication, assume risk aversion. Then $r > 0.5$ is not
1238 possible for $p = 0.5$ because then the marginal utility ratio in Eq. 4.2.2 would be at least 1 so
1239 that the right-hand side of Eq. 4.2.2 would at most be 0.5, contradiction $r > 0.5$. Applying
1240 this finding to E^c and using Eq. 2.2, $r < 0.5$ is not possible either, and $r = 0.5$ follows.

1241 Under strong risk seeking, r may differ from 0.5 for $p = 0.5$. For example, if $U(x) = e^{2.5x}$,
 1242 then $r = 0.14$ and $r = 0.86$ are optimal, and $r = 0.5$ is a local infimum, as calculations can
 1243 show. The same optimal values of r result under nonexpected utility with linear U , and with
 1244 $w(0.5) = 0.86$. Such large w -values also generate risk seeking.

1245

1246 PROOF OF THEOREM 5.2. We write π for the decision weight $W(E)$. For optimality of interior
 1247 solutions r , the first-order optimality condition for Eq. 5.4 is that

$$1248 \pi U'(a-b(1-r)^2)2b(1-r) - (1-\pi)U'(a-br^2)2br = 0,$$

1249 implying

$$1250 \pi(1-r)U'(a-b(1-r)^2) = (1-\pi)rU'(a-br^2) \quad (\text{A.1})$$

1251 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. 5.8 follows.

1252 \square

1253

1254 PROOF OF COROLLARY 6.3. Let $r > 0.5$ be optimal, and write $\pi = W(E)$. Then Eq. A.1

1255 implies

$$1256 \pi \times ((1-r)U'(a-b(1-r)^2) + rU'(a-br^2)) = rU'(a-br^2), \text{ implying}$$

$$1257 \pi = \frac{r}{r + (1-r) \frac{U'(a-b(1-r)^2)}{U'(a-br^2)}} \quad (\text{A.2})$$

1258 Applying w^{-1} to both sides yields the theorem. \square

1259

1260 In measurements of belief one first observes r , and then derives $B(E)$ from it. Corollary
 1261 6.3 gave an explicit expression. In general, it does not seem to be possible to write r as an
 1262 explicit expression of $B(E)$ or, in the case of objective probabilities with $B(E) = p$, of the
 1263 probability p .

1264

1265 PROOF OF COROLLARY 6.4. Theorem 5.2 implies that the right-hand side of Eq. 5.8 is r both
 1266 as is, and with p substituted for $B(E)$. Because Eq. 5.8 is strictly increasing in $w(B(E))$, and
 1267 w is strictly increasing too, $p = B(E)$ follows. \square

1268

1269 **Appendix B. Models for Decision under Risk and Uncertainty**

1270 For binary (two-outcome) prospects with both outcomes nonnegative, as considered in
 1271 QSRs, Eqs. 5.4 and 5.5 have appeared many times in the literature. Early references include
 1272 Allais (1953, Eq. 19.1), Edwards (1954 Figure 3), and Mosteller & Nogee (1951, p. 398).
 1273 The convenient feature that binary prospects suffice to identify utility U and the nonadditive
 1274 $w \circ B = W$ was pointed out by Ghirardato & Marinacci (2001), Gonzalez & Wu (2003), Luce
 1275 (1991, 2000), Miyamoto (1988), and Wakker & Deneffe (1996, p. 1143 and pp.1144-1145).

1276 The convenient feature that most decision theories agree on the evaluation of binary
 1277 prospects was pointed out by Miyamoto (1988), calling Eqs. 5.4 and 5.5 generic utility, and
 1278 Luce (1991), calling these equations binary rank-dependent utility. It was most clearly
 1279 analyzed by Ghirardato & Marinacci (2001), who called the equations the biseparable model.
 1280 These three works also axiomatized the model. The agreement for binary prospects was also
 1281 central in many works by Luce (e.g., Luce, 2000, Ch. 3) and in Gonzalez & Wu (2003).
 1282 Only for more than two outcomes, the theories diverge (Mosteller & Nogee 1951 p. 398;
 1283 Luce 2000, introductions to Chs. 3 and 5). Theories that also deviate for two outcomes
 1284 include betweenness models (Chew & Tan 2005), the variational model (Maccheroni,
 1285 Marinacci, & Rustichini (2006), and models with underlying multistage decompositions
 1286 (Halevy & Feltkamp 2005; Halevy & Ozdenoren 2007; Klibanoff, Marinacchi, & Mukerji
 1287 2005; Nau 2006; Olszewski 2007).

1288 We next describe some of the agreeing decision theories. Because we consider only
 1289 nonnegative outcomes, losses play no role, and we describe prospect theory only for gains
 1290 hereafter.

1291 We begin with decision under risk, with known objective probabilities $P(E)$. Expected
 1292 utility (von Neumann & Morgenstern, 1944) is the special case where w is the identity and
 1293 $B(E) = P(E)$. Kahneman & Tversky's (1979) original prospect theory, Quiggin's (1982)
 1294 rank-dependent utility, and Tversky & Kahneman's (1992) new prospect theory concern the
 1295 special case of $B(E) = P(E)$, where w now can be nonlinear. The case $B(E) = P(E)$ also
 1296 includes Gul's (1991) disappointment aversion theory.

1297 We next consider the more general case where no objective probabilities need to be
 1298 given for all events E . Expected utility is the special case where B is an additive, now
 1299 "subjective," probability and w is the identity. Choquet expected utility (Schmeidler 1989)
 1300 and cumulative prospect theory (Tversky & Kahneman 1992) start from the general

1301 weighting function W , from which B obviously results as $w^{-1}(W)$, with w the probability
 1302 weighting function for risk. The multiple priors model (Gilboa & Schmeidler 1989; Wald
 1303 1950) results with $W(E)$ the infimum value $P(E)$ over all priors P . Under Machina &
 1304 Schmeidler's (1992) probabilistic sophistication, B is an additive probability measure.
 1305

1306 **Appendix C. Experimental Instructions**

1307 This appendix will be a technical appendix on internet. It is now available at
 1308 <http://people.few.eur.nl/wakker/pdf/qsrappendixc.pdf>
 1309

1310 **References**

- 1311 Abdellaoui, Mohammed (2000), "Parameter-Free Elicitation of Utilities and Probability
 1312 Weighting Functions," *Management Science* 46, 1497–1512.
- 1313 Allais, Maurice (1953), "Le Comportement de l'Homme Rationnel devant le Risque: Critique
 1314 des Postulats et Axiomes de l'Ecole Américaine," *Econometrica* 21, 503–546.
- 1315 Allen, Franklin (1987), "Discovering Personal Probabilities when Utility Functions are
 1316 Unknown," *Management Science* 33, 542–544.
- 1317 Aragonés, Enriqueta, Itzhak Gilboa, Andrew Postlewaite, & David Schmeidler (2005), "Fact-
 1318 Free Learning," *American Economic Review* 95, 1355–1368.
- 1319 Bateman, Ian J., Alistair Munro, Bruce Rhodes, Chris Starmer, & Robert Sugden (1997),
 1320 "Does Part-Whole Bias Exist? An Experimental Investigation," *Economic Journal* 107,
 1321 322–332.
- 1322 Bernoulli, Daniel (1738), "Specimen Theoriae Novae de Mensura Sortis," *Commentarii*
 1323 *Academiae Scientiarum Imperialis Petropolitanae* 5, 175–192.
- 1324 Bleichrodt, Han & José Luis Pinto (2000), "A Parameter-Free Elicitation of the Probability
 1325 Weighting Function in Medical Decision Analysis," *Management Science* 46,
 1326 1485–1496.
- 1327 Braga, Jacinto & Chris Starmer (2005), "Preference Anomalies, Preference Elicitation, and the
 1328 Discovered Preference Hypothesis," *Environmental and Resource Economics* 32, 55–89.
- 1329 Brier, Glenn W. (1950), "Verification of Forecasts Expressed in Terms of Probability,"
 1330 *Monthly Weather Review* 78, 1–3.

- 1331 Broome, John R. (1990), "Bolker-Jeffrey Expected Utility Theory and Axiomatic
1332 Utilitarianism," *Review of Economic Studies* 57, 477–502.
- 1333 Camerer, Colin F. & Martin Weber (1992), "Recent Developments in Modelling Preferences:
1334 Uncertainty and Ambiguity," *Journal of Risk and Uncertainty* 5, 325–370.
- 1335 Charness, Gary & Dan Levin (2005), "When Optimal Choices Feel Wrong: A Laboratory
1336 Study of Bayesian Updating, Complexity, and Affect," *American Economic Review* 95,
1337 1300–1309.
- 1338 Chew, Soo Hong & Guofu Tan (2005), "The Market for Sweepstakes," *Review of Economic
1339 Studies* 72, 1009–1029.
- 1340 Clemen, Robert T. & Kenneth C. Lichtendahl (2005), "Debiasing Expert Overconfidence: A
1341 Bayesian Calibration Model," Fuqua School of Business, Duke University, Durham, NC.
- 1342 Clemen, Robert T. & Fred Rolle (2001), "In Theory ... In Practice," *Decision Analysis
1343 Newsletter* 20, No 1, 3.
- 1344 de Finetti, Bruno (1962), "Does It Make Sense to Speak of "Good Probability Appraisers"?"
1345 In Isidore J. Good (Ed.), *The Scientist Speculates: An Anthology of Partly-Baked Ideas*,
1346 William Heinemann Ltd., London.
- 1347 Dow, James & Sérgio R.C. Werlang (1992), "Uncertainty Aversion, Risk Aversion and the
1348 Optimal Choice of Portfolio," *Econometrica* 60, 197–204.
- 1349 Echternacht, Gary J. (1972), "The Use of Confidence Testing in Objective Tests," *Review of
1350 Educational Research* 42, 217–236.
- 1351 Edwards, Ward (1954), "The Theory of Decision Making," *Psychological Bulletin* 51,
1352 380–417.
- 1353 Ellsberg, Daniel (1961), "Risk, Ambiguity and the Savage Axioms," *Quarterly Journal of
1354 Economics* 75, 643–669.
- 1355 Ghirardato, Paolo & Massimo Marinacci (2001), "Risk, Ambiguity, and the Separation of
1356 Utility and Beliefs," *Mathematics of Operations Research* 26, 864–890.
- 1357 Gilboa, Itzhak (1987), "Expected Utility with Purely Subjective Non-Additive Probabilities,"
1358 *Journal of Mathematical Economics* 16, 65–88.
- 1359 Gilboa, Itzhak & David Schmeidler (1989), "Maxmin Expected Utility with a Non-Unique
1360 Prior," *Journal of Mathematical Economics* 18, 141–153.
- 1361 Gilboa, Itzhak & David Schmeidler (1999), "A *Theory of Case-Based Decisions*." Cambridge
1362 University Press, Cambridge, UK.

- 1363 Gonzalez, Richard & George Wu (1999), "On the Shape of the Probability Weighting
1364 Function," *Cognitive Psychology* 38, 129–166.
- 1365 Gonzalez, Richard & George Wu (2003), "Composition Rules in Original and Cumulative
1366 Prospect Theory," mimeo.
- 1367 Good, Isidore J. (1952), "Rational Decisions," *Journal of the Royal Statistical Society Series*
1368 *B* 14, 107–114.
- 1369 Gul, Faruk (1991), "A Theory of Disappointment Aversion," *Econometrica* 59, 667–686.
- 1370 Halevy, Yoram (2007), "Ellsberg Revisited: An Experimental Study," *Econometrica*,
1371 forthcoming.
- 1372 Halevy, Yoram & Vincent Feltkamp (2005), "A Bayesian Approach to Uncertainty
1373 Aversion," *Review of Economic Studies* 72, 449–466.
- 1374 Halevy, Yoram & Emre Ozdenoren (2007), "Uncertainty and Compound Lotteries:
1375 Calibration," working paper, University of British Columbia.
- 1376 Hansen, Lars Peter, Thomas J. Sargent, & Thomas D. Tallarini (1999), "Robust Permanent
1377 Income and Pricing," *Review of Economic Studies* 66, 873–908.
- 1378 Hanson, Robin (2002), "Wanna Bet?" *Nature* 420, November 2002, pp. 354–355.
- 1379 Harrison, Glenn W., Morten I. Lau, & M.B. Williams (2002), "Estimating Individual Discount
1380 Rates in Denmark: A Field Experiment," *American Economic Review* 92, 1606–1617.
- 1381 Hogarth, Robin M. & Hillel J. Einhorn (1990), "Venture Theory: A Model of Decision
1382 Weights," *Management Science* 36, 780–803.
- 1383 Hogarth Robin M. & Howard C. Kunreuther (1985), "Ambiguity and Insurance Decisions,"
1384 *American Economic Review, Papers and Proceedings* 75, 386–390.
- 1385 Holt, Charles A. (1986), "Preference Reversals and the Independence Axiom," *American*
1386 *Economic Review* 76, 508–513.
- 1387 Holt, Charles A. (2006), "Webgames and Strategy: Recipes for Interactive Learning," in press.
- 1388 Holt, Charles A. & Susan K. Laury (2002), "Risk Aversion and Incentive Effects," *American*
1389 *Economic Review* 92, 1644–1655.
- 1390 Huck, Steffen & Georg Weizsäcker (2002), "Do Players Correctly Estimate What Others Do?
1391 Evidence of Conservatism in Beliefs," *Journal of Economic Behavior and Organization*
1392 47, 71–85.
- 1393 Johnstone, David J. (2006), "The Value of Probability Forecast from Portfolio Theory,"
1394 School of Business, University of Sydney, Australia.

- 1395 Jouini, Elyès & Clotilde Napp (2007), “Consensus Consumer and Intertemporal Asset
1396 Pricing with Heterogeneous Beliefs,” *Review of Economic Studies* 74, 1149–1174.
- 1397 Kahneman, Daniel & Amos Tversky (1979), “Prospect Theory: An Analysis of Decision
1398 under Risk,” *Econometrica* 47, 263–291.
- 1399 Karni, Edi (2007), “A New Approach to Modeling Decision-Making under Uncertainty,
1400 *Economic Theory* 33, 225–242.
- 1401 Karni, Edi & Zvi Safra (1987), “Preference Reversal and the Observability of Preferences by
1402 Experimental Methods,” *Econometrica* 55, 675–685.
- 1403 Karni, Edi & Zvi Safra (1989), “Dynamic Consistency, Revelations in Auctions and the
1404 Structure of Preferences,” *Review of Economic Studies* 56, 421–434.
- 1405 Keren, Gideon B. (1991), “Calibration and Probability Judgments: Conceptual and
1406 Methodological Issues,” *Acta Psychologica* 77, 217–273.
- 1407 Keynes, John Maynard (1921), “*A Treatise on Probability.*” McMillan, London.
- 1408 Klibanoff, Peter, Massimo Marinacci, & Sujoy Mukerji (2005), “A Smooth Model of
1409 Decision Making under Ambiguity,” *Econometrica* 73, 1849–1892.
- 1410 Knight, Frank H. (1921), “*Risk, Uncertainty, and Profit.*” Houghton Mifflin, New York.
- 1411 Li, Wei (2007), “Changing One’s Mind when the Facts Change: Incentives of Experts and
1412 the Design of Reporting Protocols,” *Review of Economic Studies* 74,
- 1413 Luce, R. Duncan (1991), “Rank- and-Sign Dependent Linear Utility Models for Binary
1414 Gambles,” *Journal of Economic Theory* 53, 75–100.
- 1415 Luce, R. Duncan (2000), “*Utility of Gains and Losses: Measurement-Theoretical and
1416 Experimental Approaches.*” Lawrence Erlbaum Publishers, London.
- 1417 Maccheroni, Fabio, M. Marinacci, & A Rustichini (2006), “Ambiguity Aversion, Robustness,
1418 and the Variational Representation of Preferences,” *Econometrica* 74, 1447–1498.
- 1419 Machina, Mark J. (1987), “Choice under Uncertainty: Problems Solved and Unsolved,”
1420 *Journal of Economic Perspectives* 1 no 1, 121–154.
- 1421 Machina, Mark J. (2004), “Almost-Objective Uncertainty,” *Economic Theory* 24, 1–54.
- 1422 Machina, Mark J. & David Schmeidler (1992), “A More Robust Definition of Subjective
1423 Probability,” *Econometrica* 60, 745–780.
- 1424 Manski, Charles F. (2004), “Measuring Expectations,” *Econometrica* 72, 1329–1376.
- 1425 McClelland, Alastair & Fergus Bolger (1994), “The Calibration of Subjective Probabilities:
1426 Theories and Models 1980–1994.” In George Wright & Peter Ayton (eds.), *Subjective
1427 Probability*, 453–481, Wiley, New York.

- 1428 McKelvey, Richard & Talbot Page (1986), "Common Knowledge, Consensus, and Aggregate
1429 Information," *Econometrica* 54, 109–127.
- 1430 Miyamoto, J.M. (1988), "Generic Utility Theory: Measurement Foundations and Applications
1431 in Multiattribute Utility Theory," *Journal of Mathematical Psychology* 32, 357–404.
- 1432 Mosteller, Frederick & Philip Noguee (1951), "An Experimental Measurement of Utility,"
1433 *Journal of Political Economy* 59, 371–404.
- 1434 Mukerji, Sujoy (1997), "Understanding the Nonadditive Probability Decision Model,"
1435 *Economic Theory* 9, 23–46.
- 1436 Mukerji, Sujoy & Jean-Marc Tallon (2001), "Ambiguity Aversion and Incompleteness of
1437 Financial Markets," *Review of Economic Studies* 68, 883–904.
- 1438 Murphy, Allan H. & Robert L. Winkler (1974), "Subjective Probability Forecasting
1439 Experiments in Meteorology: Some Preliminary Results," *Bulletin of the American
1440 Meteorological Society* 55, 1206–1216.
- 1441 Nyarko, Yaw & Andrew Schotter (2002), "An Experimental Study of Belief Learning Using
1442 Elicited Beliefs," *Econometrica* 70, 971–1005.
- 1443 Olszewski, Wojciech (2007), "Preferences over Sets of Lotteries," *Review of Economic
1444 Studies* 74, 567–595.
- 1445 Nau, Robert F. (2006), "Uncertainty Aversion with Second-Order Utilities and Probabilities,"
1446 *Management Science* 52, 136–145.
- 1447 Palfrey, Thomas R. & Stephanie W. Wang (2007), "On Eliciting Beliefs in Strategic Games,"
1448 Division of the Humanities and Social Sciences, CalTech, Pasadena, CA 91125.
- 1449 Palmer, Tim N. & Renate Hagedorn (2006, Eds), "*Predictability of Weather and Climate.*"
1450 Cambridge University Press, Cambridge.
- 1451 Prelec, Drazen (1998), "The Probability Weighting Function," *Econometrica* 66, 497–527.
- 1452 Prelec, Drazen (2004), "A Bayesian Truth Serum for Subjective Data," *Science* 306, October
1453 2004, 462–466.
- 1454 Quiggin, John (1982), "A Theory of Anticipated Utility," *Journal of Economic Behaviour
1455 and Organization* 3, 323–343.
- 1456 Raiffa, Howard (1968), "*Decision Analysis.*" Addison-Wesley, London.
- 1457 Sandroni, Alvaro, Rann Smorodinsky, & Rakesh V. Vohra (2003), "Calibration with Many
1458 Checking Rules," *Mathematics of Operations Research* 28, 141–153.
- 1459 Savage, Leonard J. (1971), "Elicitation of Personal Probabilities and Expectations," *Journal
1460 of the American Statistical Association* 66, 783–801.

- 1461 Schmeidler, David (1989), "Subjective Probability and Expected Utility without Additivity,"
1462 *Econometrica* 57, 571–587.
- 1463 Schoemaker, Paul J.H. (1982), "The Expected Utility Model: Its Variations, Purposes,
1464 Evidence and Limitations," *Journal of Economic Literature* 20, 529–563.
- 1465 Selten, Reinhard, Abdolkarim Sadrieh, & Klaus Abbink (1999), "Money Does not Induce
1466 Risk Neutral Behavior, but Binary Lotteries Do even Worse," *Theory and Decision* 46,
1467 211–249.
- 1468 Shackle, George L.S. (1949), "A Non-Additive Measure of Uncertainty," *Review of*
1469 *Economic Studies* 17, 70–74.
- 1470 Shafer, Glenn (1976), "A *Mathematical Theory of Evidence*." Princeton University Press, NJ.
- 1471 Shiller, Robert J., Fumiko Kon-Ya, & Yoshiro Tsutsui (1996), "Why Did the Nikkei Crash?
1472 Expanding the Scope of Expectations Data Collection," *The Review of Economics and*
1473 *Statistics* 78, 156–164.
- 1474 Spiegelhalter, David J. (1986), "Probabilistic Prediction in Patient Management and Clinical
1475 Trials," *Statistics in Medicine* 5, 421–433.
- 1476 Staël von Holstein, Carl-Axel S. (1972), "Probabilistic Forecasting: An Experiment Related
1477 to the Stock Market," *Organizational Behaviour and Human Performance* 8, 139–158.
- 1478 Starmer, Chris (2000), "Developments in Non-Expected Utility Theory: The Hunt for a
1479 Descriptive Theory of Choice under Risk," *Journal of Economic Literature* 38, 332–382.
- 1480 Starmer, Chris & Robert Sugden (1991), "Does the Random-Lottery Incentive System Elicit True
1481 Preferences? An Experimental Investigation," *American Economic Review* 81, 971–978.
- 1482 Sugden, Robert (1991), "Rational Choice: A Survey of Contributions from Economics and
1483 Philosophy," *Economic Journal* 101, 751–785.
- 1484 Sugden, Robert (2004), "Alternatives to Expected Utility." In Salvador Barberà, Peter J.
1485 Hammond, & Christian Seidl, *Handbook of Utility Theory, Vol. II*, 685–755, Kluwer
1486 Academic Publishers, Dordrecht.
- 1487 Tetlock, Philip E. (2005), "*Expert Political Judgment*." Princeton University Press, NJ.
- 1488 Thaler, Richard H. & Eric J. Johnson (1990), "Gambling with the House Money and Trying
1489 to Break Even: The Effects of Prior Outcomes on Risky Choice," *Management Science*
1490 36, 643–660.
- 1491 Tversky, Amos & Daniel Kahneman (1992), "Advances in Prospect Theory: Cumulative
1492 Representation of Uncertainty," *Journal of Risk and Uncertainty* 5, 297–323.

- 1493 Tversky, Amos & Derek J. Koehler (1994), "Support Theory: A Nonextensional
1494 Representation of Subjective Probability," *Psychological Review* 101, 547–567.
- 1495 von Neumann, John & Oskar Morgenstern (1944, 1947, 1953), "*Theory of Games and*
1496 *Economic Behavior.*" Princeton University Press, Princeton NJ.
- 1497 Wakker, Peter P. (2004), "On the Composition of Risk Preference and Belief," *Psychological*
1498 *Review* 111, 236–241.
- 1499 Wakker, Peter P. & Daniel Deneffe (1996), "Eliciting von Neumann-Morgenstern Utilities
1500 when Probabilities Are Distorted or Unknown," *Management Science* 42, 1131–1150.
- 1501 Wald, Abraham (1950), "*Statistical Decision Functions.*" Wiley, New York.
- 1502 Winkler, Robert L. & Allan H. Murphy (1970), "Nonlinear Utility and the Probability
1503 Score," *Journal of Applied Meteorology* 9, 143–148.
- 1504 Wright, William F. (1988), "Empirical Comparison of Subjective Probability Elicitation
1505 Methods," *Contemporary Accounting* 5, 47–57.
- 1506 Yates, J. Frank (1990), "*Judgment and Decision Making.*" Prentice Hall, London.