# 27 Paradoxes of Rational Choice Theory

*Till Gruene-Yanoff*
University of Helsinki, Helsinki, Finland

15   **Abstract:** Rational choice theory (RCT) is beset with paradoxes. Why, then, bother with a
16   theory that raises numerous counterexamples, contradictions, and a seemingly endless stream
17   of mutually conflicting remedies? In contrast to this impression, I argue in this chapter that
18   RCT paradoxes play much more productive roles. Eight paradoxes are described in detail,
19   some of their proposed solutions are sketched out, and they are classified according to the
20   kind of paradox they pose. At their example I argue that RCT paradoxes, rather than providing
21   evidence for straightforward rejections of the theory, play important roles in education and in
22   normative research.

## Introduction

23

24   Rational choice theory (RCT) is beset with paradoxes. Why, then, bother with a theory that
25   raises numerous counterexamples, contradictions and a seemingly endless stream of mutually
26   conflicting remedies? That, at least, may be the impression of a novice student of RCT. In
27   contrast, I argue in this chapter, RCT paradoxes play a much more productive role. Rather than
28   suggesting straightforward rejections of the theory, or repellents to any newcomer, paradoxes
29   play important roles in education and in normative research.
30       RCT has a clear normative function: it offers tools for judging how people *ought* to form
31   their preferences – and by extension, how they ought to choose. A major problem is that there is
32   no hard basis against which to test normative theoretical claims – one cannot seek to falsify
33   such a theory with controlled experiments. Instead, researchers have to rely on normative
34   intuitions about assumptions and conclusions, and use theory to check whether these intui-
35   tions can be held consistently. This is where RCT paradoxes play a crucial role: they elicit
36   normative intuitions that pitch RCT assumptions and conclusions against each other. If a
37   paradox leads to a revision of the theory, it serves a research purpose. If it leads to a better
38   understanding of the assumptions and their conclusions, it serves an educational purpose.
39   Thus, many RCT paradoxes have proved, and continue, to be productive.
40       The chapter continues with a brief overview of RCT in ❯ Sect. 2, recalling the normative
41   claims it really makes. ❯ Section 3 discusses how its normative validity can be examined, and
42   the roles paradoxes play in that. ❯ Section 4 offers a classification of different kinds of
43   paradoxes. Eight selected paradoxes are surveyed in ❯ Sect. 5, sub-sectioned into paradoxes
44   of preference, belief, expected utility, and strategic interaction. Each one is explained, some of
45   its proposed solutions are sketched out, and it is classified according to the scheme proposed in
46   ❯ Sect. 4. ❯ Section 6 concludes the chapter.

## Rational Choice Theory

47

48   RCT is the dominant theoretical approach in microeconomics (although economists rarely
49   use the term "rational choice theory"). It is also widely used in other social-science disciplines,
50   in particular political science. In this context, the term rational choice theory is often associated
51   with the notion of economic "imperialism," implying that its use extends economics method-
52   ology into their fields.
53       Explicit theories of rational economic choice were first developed in the late nineteenth
54   century, and commonly linked the choice of an object to the increase in happiness an

55  additional increment of this object would bring. Early neoclassical economists (e.g., William
56  Stanley Jevons) held that agents make consumption choices so as to maximize their own
57  happiness. In contrast, twentieth-century economists disassociated RCT and the notion of
58  happiness: they presented rationality merely as maintaining a consistent ranking of alterna-
59  tives. Such a ranking is commonly interpreted as agents' desires or values.

60      Having no foundation in an ultimate end, the notion of rationality is reduced to the
61  consistent ranking of choice alternatives, the consistent derivation of this ranking from
62  evaluations of possible outcomes, and a consistency of beliefs employed in this derivation.
63  Thus, "rationality" explicated in RCT is considerably narrower and possibly sometimes at odds
64  with colloquial or philosophical notions. In philosophical contexts it often includes judgments
65  about ends, the prudent weighting of long-term versus short-term results, and insights into
66  purportedly fundamental moral principles. Nothing of this sort is invoked in RCT, which
67  simply claims that a rational person chooses actions in a manner consistent with his or her
68  beliefs and evaluations. Accordingly, a person considered "rational" in this sense may believe
69  that the moon is made of green cheese, may desire to waste his or her life, or may intend to
70  bring widespread destruction.

71      At the core of RCT is a formal framework that (1) makes the notion of preference
72  consistency precise and (2) offers formal proof that "maximizing one's utility" is identical to
73  "choosing according to a consistent preference ranking." A brief sketch of this framework
74  follows. (The framework presented here is based on von Neumann and Morgenstern (1947).
75  Alternative formal frameworks are to be found in Savage (1954) and Jeffrey (1990).)

76      Let $A = \{X_1, \ldots, X_n\}$ be a set of alternatives. Alternatives are either pure prospects or
77  lotteries. A pure prospect is a future event or state of the world that occurs with certainty.
78  For example, when purchasing a sandwich from a well-known international restaurant chain
79  I may expect certain taste experiences with near certainty. Lotteries, also called prospects under
80  risk, are probability distributions over events or states. For example, when consuming "pick-
81  your-own" mushrooms an agent faces the lottery $(X_1, p; X_2, 1-p)$, where $X_1$ denotes the
82  compound outcome (which has probability $p$) of falling ill due to poisoning and $X_2$ (with
83  probability $1-p$) the compound outcome of not doing so. More generally, a lottery $X$ consists of
84  a set of prospects $X_1, \ldots, X_n$ and assigned probabilities $p_1, \ldots, p_n$, such that $X = (X_1, p_1; \ldots X_n, p_n)$.
85  Obviously, the prospects $X_1, \ldots, X_n$ can be lotteries in themselves.

86      RCT takes preferences over actions to be evaluations of lotteries over action outcomes. Its
87  main contribution is to specify the relationship between preferences over actions, and prefer-
88  ences as well as beliefs over the compound outcomes of the respective lottery. It does so by
89  proving *representation theorems*. Such theorems show that, under certain conditions, all of an
90  agent's preferences can be represented by a numerical function, the so-called utility function.
91  Furthermore, the utility numbers of an action (i.e., lottery) $X = (X_1, p_1; \ldots X_n, p_n)$ and its
92  compound outcomes $X_1, \ldots, X_n$ are related to each other through the following principle:

$$u(X) = \sum_i p_i \times u(X_i) \tag{1}$$

93      In other words, the utility of a lottery is equal to the sum of the utilities of its compound
94  outcomes, weighted by the probability with which each outcome comes about. This is an
95  important result that significantly constrains the kind of preferences an agent can have. Of
96  course, because the representation result is a formal proof, all the constraining information
97  must already be present in the theorem's assumptions. I will sketch the main features of these

98  assumptions here. (For a detailed discussion, see the references in footnote 2. For more in-
99  depth overviews, see textbooks such as Luce and Raiffa (1957), Mas-Colell et al. (1995,
100 ❯ Chaps. 1 and ❯ 6) and Resnik (1987). Hargreaves Heap et al. (1992, pp. 3–26) give an
101 introductory treatment.)

102   RCT assumes that, at any time, there is a fixed set of alternatives $A = \{X_1, \ldots, X_n\}$ for any
103 agent. With respect to the agent's evaluation of these prospects, it assumes that agents can
104 always say that they prefer one prospect to another or are indifferent between them. More
105 specifically, it assumes that the agent has a preference ordering $\succeq$ over A, which satisfies the
106 following conditions. First, the ordering is assumed to be *complete*, that is,

$$\text{either } X_i \succeq X_j \text{ or } X_j \succeq X_i \text{ for all } X_i, X_j \in A. \tag{2}$$

107   Second, the ordering is assumed to be *transitive*, that is,

$$\text{if } X_i \succeq X_j \text{ and } X_j \succeq X_k, \text{ then also } X_i \succeq X_k \text{ for all } X_i, X_j, X_k \in A. \tag{3}$$

108   Completeness and transitivity together ensure that the agent has a so-called weak ordering
109 over all prospects.

110   The second domain in which RCT makes consistency assumptions concerns beliefs. In
111 particular, it assumes that each rational agent has a *coherent set of probabilistic beliefs*. Coher-
112 ence here means that beliefs can be represented as probability distributions that satisfy certain
113 properties. In particular, it is assumed that there is a probability function $p$ over all elements
114 of $A$, and that this function satisfies the following assumptions: first, for any $X$, $1 \geq p(X) \geq$
115 0; second, if $X$ is certain, then $p(X) = 1$; third, if two alternatives $X$ and $Y$ are mutually
116 exclusive, then $p(X \text{ or } Y) = p(X) + p(Y)$; finally, for any two alternatives $X$ and $Y$, $p(X$ and
117 $Y) = p(X) \times P(Y|X)$ – in other words the probability of the alternative "$X$ and $Y$" is identical to
118 the probability of $X$ multiplied by the probability of $Y$ *given* that $X$ is true.

119   The third domain in which rational choice theory makes consistency assumptions concerns
120 preferences over lotteries. In particular, it assumes the *independence condition*. If a prospect $X$
121 is preferred to a prospect $Y$, then a prospect that has $X$ as one compound outcome with a
122 probability $p$ is preferred to a prospect that has $Y$ as one compound with a probability $p$ and is
123 identical otherwise: that is, for all $X,Y,Z$: if $X \succeq Y$ then $(X,p;Z,1-p) \succeq (Y,p;Z,1-p)$.

124   These assumptions (together with a few others that are not relevant here) imply that
125 preferences over lottery prospects $X = (X_1,p_1;\ldots X_n,p_n)$ are represented by a utility function
126 such that for all $X,Y$:

$$X \succeq Y \Leftrightarrow \sum_i [p_i \times u(X_i)] \geq \sum_i [p_i \times u(Y_i)]. \tag{4}$$

127   This formal result has been given different interpretations. My focus in the following is on
128 the *normative* interpretation of RCT.

## Normative Validity and the Role of Paradox in RCT
129

130 RCT is often interpreted as a theory of how people *ought* to form their preferences – and by
131 extension, how they ought to choose (for a history of this approach, see Guala 2000). Although
132 the normative content of the theory is limited to the norms of a consistent ranking of choice

133 alternatives, I showed in the previous section that this notion of consistency depends on
134 a number of substantial axioms. This raises the question of the *normative validity* of these
135 axioms: why ought people to choose in accordance with them?

136  Various attempts have been made to defend the normative validity of RCT and its axioms.
137 The most prominent justifications are pragmatic: they seek to show that agents who fail to
138 retain RCT-consistency will incur certain losses. Two well-known examples are the *money*
139 *pump* and the *Dutch book* arguments (for more on this and other normative justifications, see
140 Hansson and Grüne-Yanoff 2009, ❯ Sect. 1).

141  Interpreted literally, neither the money pump nor the Dutch book is very convincing.
142 An agent could simply refuse to accept money-pumping trades or Dutch-booking bets. Thus,
143 rationality does not literally require one to be willing to wager in accordance with RCT.
144 Defenders of pragmatic justifications may argue that money pumps or Dutch books reveal
145 possible vulnerability from RCT-violations: a RCT violator might have an incentive to accept
146 a money pump trade or a Dutch book bet, while a RCT-abider does not. However, even such
147 a hypothetical interpretation is problematic. For example, one could deny that the situations
148 considered are normatively relevant to actual preferences. Consequently, it could be argued
149 that norms of preference consistency are primitive in the sense that they are not derived from
150 anything, and in particular not from pragmatic considerations.

151  Instead, some argue that normative judgments arise directly through human intuition,
152 guided by reflection. These judgments are grounded in characteristic human responses of
153 an emotional or motivating kind. (Such a view does not presuppose a non-cognitivist account
154 of normative judgment. At least on the epistemological level, even cognitivist theories of
155 normativity are likely to appeal to something like natural human responses – no doubt refined
156 by education and reason – to explain how we identify moral facts and evaluate moral claims.)
157 While considerations such as the money pump or the Dutch book may elicit such intuitions, it
158 would be misguided to assume that pragmatic considerations form their basis. Rather, nor-
159 mative intuitions themselves are basic, and form the basis of normative validity judgments of
160 RCT, in this view.

161  So much the worse for the normative validity of RCT axioms, one might be tempted to
162 reply. To be sure, our emotional or motivating responses to questions of preference consistency
163 often differ and are contradictory. Hence, it seems to follow that any proposed set of axioms is
164 nothing more than the expression of a subjective intuition, fuelled at best by positional or
165 rhetorical power.

166  Defenders of a stronger validity claim may respond in at least two ways to this challenge.
167 First, they may point out that normative intuitions are not merely claimed to be valid
168 individually, but rather that RCT makes a claim about the normative validity of the whole
169 set of assumptions *and all the results deduced from it.* For example, the maximization of
170 expected utility is a consequence of standard RCT axioms, not an axiom itself. If one has
171 doubts about the normative validity of this conclusion, one has to trace it back to these axioms,
172 re-check their validity, and weight one's doubts in the conclusion against one's confidence in
173 them. This view of normativity thus rests on the idea of a reflective equilibrium: we "test"
174 various parts of our system of normative intuitions against the other intuitions we have made,
175 looking for ways in which some of them support others, seeking coherence among the widest
176 set, and revising and refining them at all levels when challenges to some arise from others (for
177 more on the method of reflective equilibrium, see Daniels 2008).

Second, the defender may point out that normative intuitions are widely accepted only if they withstand being tested in a communally shared effort of "normative falsification" (Guala 2000). Savage described this effort as follows:

▶ In general, a person who has tentatively accepted a normative theory must conscientiously study situations in which the theory seems to lead him astray; he must decide for each by reflection – deduction will typically be of little relevance – whether to retain his initial impression of the situation or to accept the implications of the theory for it. (Savage 1954, p. 102)

Theorists have to engage in thought experiments in order to elicit these normative intuitions – or "initial impressions," as Savage calls them. Thereby they investigate their normative intuitions in as wide a scope of hypothetical situations as possible, either challenging or confirming particular normative judgments. At the end of this process they only use the intuitions that hold up against normative falsification to challenge the theory:

▶ If, after thorough deliberation, anyone maintains a pair of distinct preferences that are in conflict with the sure-thing principle, he must abandon, or modify, the principle; for that kind of discrepancy seems intolerable in a normative theory. (Savage 1954, p. 102)

Normative falsification and reflective equilibrium thus go hand in hand: the former generates "corroborated" normative intuitions, and the latter weighs the importance of these intuitions against conflicting intuitions in the theory under scrutiny.

How, then, does one go about normative falsification? How are "situations" constructed in which one obtains "initial impressions" that conflict with the theory? This is where RCT paradoxes come into play. These paradoxes are exemplar narratives of situations that have posed problems for RCT, many of which have been discussed amongst experts for decades. Sometimes agreed-upon solutions exist, and the paradox is used only for pedagogical purposes – to increase understanding of the theory or to illustrate the process of thought experimentation. At other times, competing solutions are offered, some of which may threaten the current theory. In that case, RCT paradoxes constitute the laboratory equipment of ongoing decision-theoretical research.

## The Notion of Paradox

Philosophers have distinguished between two accounts of paradoxes. The *argumentative model*, proposed by Quine (1966) and Sainsbury (1988), defines a paradox as an argument that appears to lead from a seemingly true statement or group of statements to an apparent or real contradiction, or to a conclusion that defies intuition. To resolve a paradox, on this account, is to show either (1) that the conclusion, despite appearances, is true, that the argument is fallacious, or that some of the premises are false, or (2) to explain away the deceptive appearances. The *non-argumentative model*, proposed by Lycan (2010), defines a paradox as an inconsistent set of propositions, each of which is very plausible. To resolve a paradox under this account is to decide on some principled grounds which of the propositions to abandon.

Consequently, the argumentative model allows distinguishing different kinds of paradoxes. Quine divides them into three groups. A *veridical paradox* produces a conclusion that is valid, although it appears absurd. (Quine thought of paradoxes pertaining to the truth of deductive

219 statements. In contrast, the validity of decision-theoretic assumptions and conclusions con-
220 cerns normative validity, which may or may not be reducible to truth. I will therefore use
221 "validity" where Quine spoke of "truth.") For example, the paradox of Frederic's birthday in
222 *The Pirates of Penzance* establishes the surprising fact that a 21-year-old would have had only
223 five birthdays had he been born in a leap year on February 29.

224 A *falsidical paradox* establishes a result that is actually invalid due to a fallacy in the
225 demonstration. DeMorgan's invalid mathematical proof that 1 = 2 is a classic example, relying
226 on a hidden division by zero.

227 Quine's distinction here is not fine enough for the current purposes. A falsidical paradox
228 in his terminology, so it seems, can be the result of two very different processes. A genuine
229 falsidical paradox, I suggest, identifies the root of the invalidity of the conclusion in the
230 invalidity of one or more of the assumptions. In contrast, what I call an *apparent paradox*
231 establishes the root of the invalidity of the conclusion in the unsoundness of the argument.

232 A paradox that is in neither class may be an *antinomy*, which reaches a self-contradictory
233 result by properly applying accepted ways of reasoning. Antinomies resist resolutions: the
234 appearances cannot be explained away, nor can the conclusion be shown to be valid, some
235 premises shown to be invalid, or the argument shown to be unsound. Antinomies, Quine says,
236 "bring on the crisis in thought" (1966, p. 5). They show the need for drastic revision in our
237 customary way of looking at things.

238 The non-argumentative account rejects Quine's classification, pointing out his assump-
239 tion of an intrinsic direction in the relationship between "assumptions" and "conclusions."
240 This, so Lycan argues, may give the wrong impression that certain kinds of paradoxes are to
241 be solved in particular ways. Conversely, he points out that two theorists may disagree on
242 whether a paradox is veridical or not: "one theorist may find the argument veridical while
243 the other finds the 'conclusion's' denial more plausible than one of the 'premises'" (Lycan   Au1
244 2010, p. 3). In what follows, I will make use of Quine's classification. Nevertheless, I stress –
245 in agreement with Lycan – that it is only to be understood as an indicator of how the
246 majority of theorists have sought resolution, not as a claim about the intrinsic nature of the
247 paradox itself.

## The Paradoxes

248

249 Below I survey a selection of paradoxes that are currently relevant to RCT. By relevant here
250 I mean that they challenge one or more of the RCT axioms that are currently in wide use. For
251 more comprehensive literature on paradoxes, also in RCT, see Richmond and Sowden (1985),
252 Diekmann and Mitter (1986), Sainsbury (1988), and Koons (1992).

253 The survey is structured according to the aspect of RCT under challenge. As it turns out, it is
254 not always clear which axiom is being challenged. I therefore divide the subsections into
255 paradoxes of preferences, belief, expected utility maximization, and strategic choice.

### Preferences

256

257 Of the many paradoxes challenging assumptions about preferences I will survey two: the Sorites
258 Paradox applied to preference transitivity and Allais' paradox.

259    Sorites paradoxes are arguments that arise from the indeterminacy surrounding the limits
260    of application of the predicates involved (for a general overview, see Hyde 2008). The Sorites
261    scheme has been applied to RCT in order to cast doubt on the rationality of the transitivity of
262    preference. Quinn's (1990) version goes as follows:

263    A person (call him *the self-torturer*) is strapped to a conveniently portable machine, which
264    administers a continuous electric current. The device has 1,001 settings: 0 (off) and 1 . . . 1,000,
265    of increasing current. The increments in current are so tiny that he cannot feel them. The self-
266    torturer has time to experiment with the device so that he knows what each of the settings feels
267    like. Then, at any time, he has two options: to stay put or to advance the dial one setting.
268    However, he may advance only one step each week, and he may never retreat. At each advance
269    he gets $10,000.

270    Since the self-torturer cannot feel any difference in comfort between adjacent settings, he
271    appears to have a clear and repeatable reason to increase the current each week. The trouble is
272    that there are noticeable differences in comfort between settings that are sufficiently far apart.
273    Eventually, he will reach settings that will be so painful that he would gladly relinquish his
274    monetary rewards and return to zero.

275    The paradox lies in the conclusion that the self-torturer's preferences are intransitive.
276    All things considered, he prefers 1–0, 2–1, 3–2, and so on, but certainly not 1,000–1. Further-
277    more, there seems to be nothing irrational about these preferences.

278    ▶   The self-torturer's intransitive preferences seem perfectly natural and appropriate given his
279        circumstances. (Quinn 1990, p. 80)

280    If this were correct, the normative validity of the transitivity axiom would be in doubt.
281    Defenders of transitivity argue that there is a mistake either in the conception of the decision
282    situation or in the process of evaluation that leads to the intransitive preferences. Both
283    Arntzenius and McCarthy (1997) and Voorhoeve and Binmore (2006) follow the first option
284    in rejecting the implicit assumption that there is a "least-noticeable difference": a magnitude of
285    physical change so small that human beings always fail to detect a difference between situations
286    in which a change smaller than this magnitude has or has not occurred.

287    Instead, they argue that it is rational for the self-torturer to take differences in long-run
288    frequencies of pain reports into account. In other words, when repeatedly experimenting with
289    the machine he may well experience different amounts of pain at the same notch. He will
290    represent this information about how a notch feels by means of a distribution over different
291    levels of pain: two notches "feel the same" only if they have the same distribution. Then it is
292    implausible that all adjacent notches feel the same when the self-torturer runs through them in
293    ascending order, and the intransitivity disappears.

294    Thus, Quinn's paradox has been treated as an apparent paradox: an implicit, illegitimate
295    assumption of the derivation – of a "least-noticeable difference" – is exposed, and the
296    dependency of the deductive conclusion on this assumption is shown.

297    *Allais' Paradox* (Allais 1953) sets up two specific choices between lotteries in order to
298    challenge the sure-thing principle (an axiom in Savage's decision theory, related to the axiom of
299    independence). This choice experiment is described in ❯ *Table 27.1*.

300    RCT prescribes that agents choose *C* if they have chosen *A* (and vice versa), and that they
301    choose *D* if they have chosen *B* (and vice versa). To see this, simply re-partition the prizes of the
302    two problems as follows: Instead of "2,400 with certainty" in *B*, partition the outcome such that
303    it reads "2,400 with probability 0.66" and "2,400 with probability 0.34." Instead of "0 with

t1.1 ◘ **Table 27.1**

**Allais' two pairs of choices**

| Choice problem 1 – choose between: | | | | | |
|---|---|---|---|---|---|
| A: | $2,500 | With probability 0.33 | B: | $2,400 | With certainty |
| | $2,400 | With probability 0.66 | | | |
| | $0 | With probability 0.01 | | | |
| Choice problem 2 – choose between: | | | | | |
| C: | $2,500 | With probability 0.33 | D: | $2,400 | With probability 0.34 |
| | $0 | With probability 0.67 | | $0 | With probability 0.66 |

t1.2
t1.3
t1.4
t1.5
t1.6
t1.7
t1.8

t2.1 ◘ **Table 27.2**

**The redescribed choice pairs**

| Choice problem 1* – choose between: | | | | | |
|---|---|---|---|---|---|
| A: | $2,500 | With probability 0.33 | B*: | $2,400 | With probability 0.34 |
| | $0 | With probability 0.01 | | | |
| Choice problem 2* – choose between: | | | | | |
| C*: | $2,500 | With probability 0.33 | D: | $2,400 | With probability 0.34 |
| | $0 | With probability 0.01 | | | |

t2.2
t2.3
t2.4
t2.5
t2.6
t2.7

304 probability 0.67" in *C*, partition the outcome such that it reads "0 with probability 0.66" and
305 "0 with probability 0.01." Of course, these are just redescriptions that do not change the nature
306 of the choice problem. They are shown in ❷ *Table 27.2*.
307   Through this redescription, we now have an outcome "2,400 with probability 0.66" in
308 both *A* and *B**, and an outcome "0 with probability 0.66" in both *C** and *D*. According to the
309 RCT independence condition, these identical outcomes can be disregarded in the deliberation,
310 but once they are disregarded it becomes clear that option *A* is identical to option *C**　and
311 option *B** is identical to option *D*. Hence, anyone choosing *A* should also choose *C* and anyone
312 choosing *B* should also choose *D*.
313   This result has been found both empirically and normatively challenging for RCT. (In sharp
314 contrast to the RCT result, in an experiment involving 72 people, 82% of the sample chose *B*
315 and 83% chose *C* (Kahneman and Tversky 1979).) On the normative level, many people seem
316 to have intuitions contradicting the above conclusions:

317 ► When the two situations were first presented, I immediately expressed preference for Gamble 1
318   [*A*] as opposed to Gamble 2 [*B*] and for Gamble 4 [*D*] as opposed to Gamble 3 [*C*], and I still feel an
319   intuitive attraction to these preferences. (Savage 1954, p. 103)

320   This empirical and normative challenge has given rise to a number of alternatives,
321 including prospect theory (Kahneman and Tversky 1979), weighted utility (Chew 1983) and

322 rank-dependent expected utility (Quiggin 1982). However, the normative validity of these
323 theories is often controversial. Many decision theorists have rather followed Savage, who
324 despite his initial intuitions decided that the sure-thing axiom was, after all, correct. He arrived
325 at this by redescribing Allais' choice situation in yet another form, observing a change of
326 preference from $C$ to $D$ in this case, and concluded that "in revising my preferences between
327 Gambles 3 [$C$] and 4 [$D$] I have corrected an error" (Savage 1954, p. 103). Decision theorists
328 following Savage have thus treated Allais' paradox as a veridical paradox: the initial impression
329 that the conclusion is absurd is explained away, and the theoretical conclusion is confirmed
330 to be correct.

## Belief

331

332 I discuss only one paradox of belief here, namely the *Monty Hall problem*. It is posed as follows:

333 ▶   Suppose you're on a game show, and you're given the choice of three doors: Behind one door is
334     a car; behind the others, goats. You pick a door, say No. 1, and the host, who knows what's behind
335     the doors, opens another door, say No. 3, which has a goat. He then says to you, ''Do you want to
336     pick door No. 2?'' Is it to your advantage to switch your choice? (vos Savant 1990).

337     After picking a door at random, it may seem that it is rational to believe that the remaining
338 door holds the car with probability $\frac{1}{2}$. After all, either the chosen door or the other one
339 conceals the prize – so should both doors rather be assigned an equal probability of holding
340 the car?

341     They should not. At first, before the contestant picks a door, it is rational for him or her to
342 believe that the car is behind any of them with probability $\frac{1}{3}$, knowing that the host will be able
343 to open a door not holding the prize since at least one of the other doors must conceal a goat.
344 Therefore, when the host opens a door the contestant does not learn anything relevant to his or
345 her belief in having chosen the winning door – it remains at $\frac{1}{3}$. Now the offered swap is
346 equivalent to the opportunity of opening both other doors – and he or she should rationally
347 believe that this offers a $\frac{2}{3}$ probability of winning the car. Hence, it is advantageous to swap.

348     The Monty Hall problem clearly falls into the class of veridical paradoxes: the argument is
349 correct and does not rely on implicit illegitimate assumptions, and the conclusion, despite
350 appearances, is valid. Furthermore (unlike in most other paradoxes discussed here), the result
351 can be experimentally confirmed: the frequency of the prize being behind the other door is
352 observed indeed to converge to $\frac{2}{3}$ (hence there is even a pragmatic confirmation of the
353 normative intuition). The striking thing about this paradox is that the presentation of the
354 correct answer initially created a huge outcry, not least from academically trained mathe-
355 maticians and logicians (see http://www.marilynvossavant.com/articles/gameshow.html for a
356 selection). The correct solution appeared to be false, but this appearance was explained away
357 with the help of standard probability theory.

## Expected Utility

358

359 I will now discuss three paradoxes that challenge some fundamental assumptions about the
360 rationality of expected utility maximization: the Ellsberg paradox, Newcomb's problem and the
361 Envelope paradox.

362  The *Ellsberg paradox* (Ellsberg 1961) goes as follows. An urn contains 30 red balls and 60
363  other balls that are either black or yellow. You do not know how many black or yellow balls
364  there are, but you do know that the total number of black balls plus the total number of yellow
365  balls equals 60. The balls are well mixed so that each individual one is as likely to be drawn as
366  any other. You are now given a choice between two gambles ( ❱ *Table 27.3*).                    Au2
367  You are also given the choice between these two gambles with regard to a different draw
368  from the same urn ( ❱ *Table 27.4*).
369  Standard (Bayesian) decision theory postulates that when choosing between these gambles,
370  people assume a probability that the non-red balls are yellow versus black, and then compute
371  the expected utility of the two gambles. This leads to the following line of reasoning. The prizes
372  are exactly the same. Hence, according to expected utility theory, a rational agent (weakly)
373  prefers *A* to *B* if and only if he or she believes that drawing a red ball is at least as likely as
374  drawing a black ball. Similarly, a rational agent (weakly) prefers *C* to *D* if and only if he or she
375  believes that drawing a red or yellow ball is at least as likely as drawing a black or yellow ball.
376  Now, if drawing a red ball is at least as likely as drawing a black ball, then drawing a red or
377  yellow ball is also at least as likely as drawing a black or yellow ball. Thus, supposing that
378  a rational agent (weakly) prefers *A* to *B*, it follows that he or she will also (weakly) prefer *C* to *D*,
379  whereas supposing instead a weak preference for *D* over *C*, it follows that he or she will also
380  weakly prefer *B* to *A*. When surveyed, however, most people (strictly) prefer Gamble A to
381  Gamble B and Gamble D to Gamble C. Furthermore, they often insist on this choice, even if the
382  theory is explained to them. Therefore, the normative validity of some assumptions of RCT
383  seems in question.
384  Ellsberg's paradox poses an interesting challenge to RCT. On the one hand, some scholars
385  have insisted that the standard solution is correct, making it a veridical paradox whose
386  paradoxical impression is explained away. Fox and Tversky (1995), for example, offer an
387  empirical explanation of why people may be biased in their decision-making through an
388  impression of comparative ignorance. Such bias, of course, has no normative validity: it
389  only explains why people have wrong intuitions about what should be chosen. On the other
390  hand, others have argued that this ambiguity aversion is part of a rational decision, in a similar
391  way as a risk aversion is (Schmeidler 1989). Such a position would suggest that the Ellsberg
392  paradox is falsidical, brought about by assuming away the (rationally relevant) impact of
393  ambiguity aversion.

t3.1  ◼ **Table 27.3**
**Ellsberg's first pair of choices**

t3.2

| Gamble *A* | Gamble *B* |
|---|---|
| You receive $100 if you draw a red ball | You receive $100 if you draw a black ball |

t3.3

t4.1  ◼ **Table 27.4**
**Ellsberg's pair of choices**

t4.2

| Gamble *C* | Gamble *D* |
|---|---|
| You receive $100 if you draw a red or yellow ball | You receive $100 if you draw a black or yellow ball |

t4.3

394    *Newcomb's Problem* (Nozick 1969) involves an agent's choosing either an opaque box
395    or the opaque and a transparent box. The transparent box contains one thousand dollars
396    ($T$) that the agent plainly sees. The opaque box contains either nothing ($0) or one million
397    dollars ($M$), depending on a prediction already made concerning the agent's choice. If the
398    prediction was that the agent would take both boxes, then the opaque box will be empty, and if
399    it was that the agent would take just the opaque box then the opaque box would contain
400    a million dollars. The prediction is reliable. The agent knows all these features of the decision
401    problem.
402    ❯ *Table 27.5* displays the agent's choices and their outcomes. A row represents an option,
403    a column a state of the world, and a cell an option outcome in a state of the world.
404    Standard RCT posits that a rational agent should choose the option that maximizes
405    expected utility. This approach recommends taking only one box, for the following reasons.
406    First, the prediction is reliable. In other words, if the agent chooses only one box, then the
407    probability that "take one box" was predicted is high. Similarly, if the agent chooses two boxes,
408    then the probability that "take two boxes" was predicted is high. Hence the probability of
409    outcome $M given that the agent chose only one box will be high, and the probability of
410    outcome $M + $T given that the agent chose two boxes will be low – sufficiently low in most
411    plausible cases for the expected utility of "taking one box" to be higher that that of "taking two
412    boxes." Hence one-boxing is the rational choice according to the principle of expected-utility
413    maximization.
414    Yet this recommendation violates two deeply entrenched intuitions. First, it violates
415    the principle of dominance, according to which an agent prefers one action to another if he
416    or she prefers every outcome of the first action to the corresponding outcomes of the second.
417    The normative validity of dominance is widely agreed upon. Yet, two-boxing clearly dominates
418    one-boxing in this sense. Consequently, RCT violates dominance.
419    Second, it violates the intuition that actions should be chosen on the basis of their causal
420    effects rather than their probabilistic correlations to benefits or drawbacks. Because the
421    prediction is made before the agent chooses, the choice has no causal impact on it, and the
422    probabilistic correlation should not matter.
423    This analysis has motivated a reformulation of decision theory on causal rather than
424    evidential grounds. In various accounts, causal decision theorists seek to represent causal
425    influence with their probability functions rather than with mere probabilistic correlation
426    (see e.g., Gibbard and Harper 1981; Skyrms 1980; Joyce 1999). They clearly see Newcomb's
427    problem as a falsidical paradox based on the misspecification of a decision maker's
428    relevant beliefs.
429    Some authors opposing the causal approach hold that it yields the wrong choice in
430    Newcomb's problem, in other words two-boxing rather than one-boxing. Horgan (1985)

t5.1    ◾ **Table 27.5**

**Newcomb's problem**

| | Prediction of one-boxing | Prediction of two-boxing |
|---|---|---|
| Take one box | $M | $0 |
| Take two boxes | $M + $T | $T |

and Horwich (1987), for example, argue that one-boxers fare better than two-boxers, and that one-boxing is therefore the rational choice of action. They both see Newcomb's problem as a veridical paradox, and propose explaining away the conflicting intuition about dominance and causal influence with reference to pragmatic success. Causal decision theorists, in turn, reject the relevance of these pragmatic considerations for the validity of the above intuitions. They insist that Newcomb's problem is an unusual case, which rewards irrationality: one-boxing is irrational even if one-boxers prosper.

The *two-envelope paradox* goes as follows. You are asked to make a choice between two envelopes. You know that one of them contains twice the amount of money as the other, but you do not know which one. You arbitrarily choose one envelope – call it Envelope *A* – but do not open it. Call the amount of money in that envelope *X*. Since your choice was arbitrary, there is a 50–50 chance that the other envelope (Envelope *B*) will contain more money, and a 50–50 chance that it will contain less. Would you now wish to switch envelopes?

Calculating the apparent expected value of switching proceeds as follows. Switching to *B* will give you a 50% chance of doubling your money and a 50% chance of halving it. Thus it seems that the expected value of switching to *B* is $E(Y- X) = 0.5*1/2X + 0.5*2X - X = 0.25X$. Hence, switching to *B* will give you a 25% higher expected return than sticking with *A*. This seems absurd. First, many people have an intuition that one should be indifferent between *A* and *B* as long as the envelope remains unopened. Second, once you have switched to *B* in line with the above argument, a symmetrical calculation could persuade you to switch back to *A*. Therein lies the paradox.

It is now widely agreed that the expected gain from switching, $E(Y-X)$, is mathematically undefined because the value of the infinite sum of all probability-weighted values of *Y-X* depends on the order of summation (Meacham and Weisberg 2003). However, the conclusions from this observation differ widely. Clark and Shackel (2000) argue that there is a "correct" order of summation, which results in a zero infinite sum, and that this result justifies indifference before opening the envelope. In contrast, Meacham and Weisberg (2003) express reservations about selecting the "correct" order of summation: because the expected gain from switching is undefined, standard decision theory does not rank switching against keeping.

Dietrich and List (2005) go along a different route and offer an axiomatic justification for indifference before opening without appeal to infinite expectations. They supplement standard decision theory with an additional axiom, the "indifference principle," according to which if two lotteries have identical distributions, a rational agent is indifferent between them. From this they are able to deduce a justification for indifference before opening. All three of these responses, although formulated against each other, consider the two-envelope paradox falsidical: Clark and Shackel and Dietrich and List introduce additional assumptions, which yield the intuitive conclusion, whereas Mechaem and Weisberg insist that the argument is fallacious, and no conclusion is warranted from the given assumptions.

## Strategic Interaction

Game theory is closely related to RCT. Although it requires certain assumptions beyond those of the standard RCT axioms, its models give additional significance to those standard RCT axioms that also play a role in game theory. For this reason, I include two game-theoretic paradoxes here: the Prisoners' dilemma and the paradox of common knowledge.

474    The One-Shot *Prisoners' Dilemma* has attracted much attention because standard game-
475 theoretic solution concepts unanimously advise each player to choose a strategy that will result
476 in a Pareto-dominated outcome. It goes as follows.

477    Two gangsters who have been arrested for robbery are placed in separate cells. Both care
478 much more about their personal freedom than about the welfare of their accomplice.
479 A prosecutor offers each the following deal: choose to confess or remain silent. If one prisoner
480 confesses and the accomplice remains silent all charges against the former are dropped
481 (resulting in a utility of 3 in ❯ *Table 27.6* – the first number in each cell is the utility of this
482 outcome for the player who chooses between rows, and the second number, the utility for the
483 player who chooses between columns). If the accomplice confesses and the first prisoner
484 remains silent however, the latter will do time (utility 0 in ❯ *Table 27.6*). If both confess,
485 each will get reduced sentences (utility 1), and if both remain silent the prosecutor has to settle
486 for token sentences for firearms possession (utility 2) (for extensive discussion and a literature
487 review, see Kuhn 2009).

488    The choice situation is solved by appeal to a simple dominance argument. For each
489 player, if the other player stays silent it is better to confess than to stay silent. If the other
490 player confesses, it is also better to confess than to stay silent. Hence, no matter what the other
491 player does, it is always better to confess.

492    This result is often described as paradoxical in the following sense. The outcome obtained
493 when both confess, although it is rational for each to do so, is worse for each than the outcome
494 they would have obtained had both remained silent. Both would prefer to reach the outcome
495 "stay silent, stay silent," but their individually rational actions led them to the inferior result
496 "confess, confess." (To add some more urgency to this example, consider the structurally
497 similar problem of the "tragedy of the commons," according to which multiple individuals
498 acting independently and rationally will ultimately deplete a shared limited resource even when
499 it is clear that it is not in anyone's long-term interest for this to happen (Hardin 1968).) How
500 can such an inferior outcome be the result of rational decisions?

501    Some authors argue that the Prisoners' Dilemma indeed exposes a limitation of RCT
502 rationality. Gauthier (1986), for example, suggests that, instead of always confessing, it would
503 be rational for players to commit to playing cooperatively when faced with other cooperators
504 who are equally committed to not exploiting one another's good will. This argument crucially
505 depends on player confidence in that most players are clearly identifiable as being committed
506 to cooperating or not. Whether such a belief can be rationally justified is questionable, and
507 with it the whole solution to the dilemma.

508    The majority of authors see no conceptual problem in the Prisoners' Dilemma. The
509 assumptions say nothing about the necessary selfishness of the players (charity organizations
510 may also find themselves in such situations!), and no other illegitimate assumptions are

t6.1    **◻ Table 27.6**

**The Prisoners' dilemma**

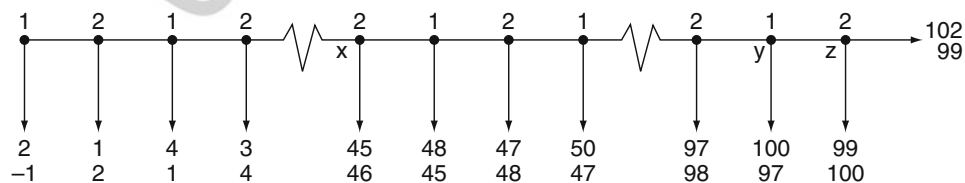| | Stay silent | Confess |
|---|---|---|
| *Stay silent* | 2,2 | 0,3 |
| *Confess* | 3,0 | 1,1 |

511 evident. The argument itself is valid, and the conclusion is not contradicted by normative
512 intuitions. The only problem is terminological: some people chafe against the idea that the
513 conclusion is supposed to be the result of rational decision-making. However, they simply
514 subscribe to a different concept of rationality that RCT does not support. Hence the
515 Prisoners' Dilemma is a veridical paradox, whose paradoxical nature relies on terminological
516 ambiguity.

517 *Backward induction* is the process of reasoning backward in time, from the end of
518 a problem, to determine a sequence of optimal actions. It proceeds by first considering the
519 latest time a decision can be made and choosing what to do in any situation at that time. One
520 can then use this information to determine what to do at the second-to-last time for the
521 decision. This process continues backward until one has determined the best action for every
522 possible situation at any point in time.

523 Let us take a concrete example, the "centipede game." This game progresses from left
524 to right in ❯ *Fig. 27.1*. Player 1 (female) starts at the extreme left node, choosing to end the
525 game by playing *down*, or to continue (giving player 2, male, the choice) by playing *right*. The
526 payoffs are such that at each node it is best for the player whose move it is to stop the game if
527 and only if he or she expects it to end at the next stage if he or she continues (if the other player
528 stops the game or if it is terminated). The two zigzags stand for the continuation of the payoffs
529 along those lines. Now backward induction advises resolving the game by starting at the last
530 node *z*, asking what player 2 would have done had he ended up there. A comparison of player
531 2's payoffs for the two choices implies that he should have rationally chosen *down*. Given
532 common knowledge of rationality, the payoffs that result from this choice of *down* can be
533 substituted for node *z*. Let us now move backwards to player 1's decision node. What would she
534 have done had she ended up at node *y*? Given player 2's choice of *down*, she would have chosen
535 *down*, too. This line of argument then continues all the way back to the first node. Backward
536 induction thus recommends player 1 to play *down* at the first node.

537 What, then, should player 2 do if he actually found himself at node *x*? Backward induction
538 tells him to play "down," but backward induction also tells him that if player 1 were rational he
539 would not be facing the choice at node *x* in the first place. This is not a problem in that
540 Backward induction predicts that player 2 will never find himself at *x* unless both players
541 are irrational.

542 Yet what does this imply for the Backward-induction reasoning process itself? Backward
543 induction requires the players to *counterfactually* consider out-of-equilibrium play. For exam-
544 ple, player 1, according to Backward induction, should choose *down* at node 1, because she
545 knows that player 2 would have chosen *down* at node 2, which in turn she knows because she
546 would have chosen *down* at node 3, and so on, because ultimately she knows that player 2
547 would have chosen *down* at *z*. She knows this because she knows that player 2 is rational,



Au3   ◼ **Fig. 25.1**

548 and that player 2 knows that she is rational, and so on. In the language of game theory,
549 because rationality is *common knowledge* amongst the players, backward induction applies
550 (for more on common knowledge, including a formal treatment of this paradox, see
551 Vanderschraaf and Sillari 2009).

552     Given common knowledge of rationality, each player can affirm the counterfactual "A
553 rational player finding himself or herself at any node in the centipede would choose *down*." Yet
554 we also concluded that if a player *finds* himself or herself at any node with an index number
555 larger than two then both player and opponent know that they *are* not rational. What if that
556 conclusion also made true the counterfactual "if a player *found* himself or herself at any node
557 with an index larger than two then both player and opponent *would* know that they are not
558 rational"? If it did, Backward induction would break down: it requires (1) common knowledge
559 and (2) counterfactual consideration of how players would choose if they found themselves at
560 nodes with indices larger than two. However, (2) implies that both player and opponent would
561 know that they are not rational, contradicting (1). Herein lies the paradox (Pettit and Sugden
562 1989; Bicchieri 1989).

563     This is an intensely discussed problem in game theory and philosophy. There is space here
564 only to sketch two possible solutions. According to the first, common knowledge of rationality
565 implies backward induction in games of perfect information (Aumann 1995). This position
566 is correct in that it denies the connection between the indicative and the counterfactual
567 conditional. Players have common knowledge of rationality, and they are not going to lose it
568 regardless of the counterfactual considerations they engage in. Only if common knowledge
569 were not immune to evidence, and would be revised in the light of the opponents' moves,
570 might this sufficient condition for backward induction run into the *conceptual problem*
571 sketched above. However, common knowledge, by definition, is not revisable, and thus the
572 argument has to assume a *common belief* in rationality instead. If one looks more closely at the
573 versions of the above argument (e.g., Pettit and Sugden 1989) it becomes clear that they employ
574 the notion of common belief rather than common knowledge. Hence, the backward-induction
575 paradox is only apparent: the argument that led to the seemingly contradictory conclusion
576 is unsound.

577     The second potential solution obtains when one shows, as Bicchieri (1993, ❯ Chap. 4)
578 does, that limited knowledge (and *not* common knowledge) of rationality and of the structure
579 of the game suffice for backward induction. All that is needed is that a player at each
580 information set knows what the next player to move knows. This condition does not get
581 entangled in internal inconsistency, and backward induction is justifiable without conceptual
582 problems. In that case, the backward-induction paradox is falsidical.

## Further Research

583

584 I have surveyed eight paradoxes of RCT. There is considerable divergence among them, under
585 a rather rough classificatory scheme, even in this small selection. First, there are the veridical
586 paradoxes, like the Prisoner's dilemma, the paradoxical nature of which rests merely on
587 terminological ambiguity. Ways of explaining away the paradoxical appearance of other
588 veridical paradoxes such as the Monty Hall problem are obvious, but are baffling to the novice.
589 They can still serve an educational purpose, however, in that studying them clarifies the
590 meaning of the assumptions and the derivation of the conclusion.

591  Second, there are (relatively) clear cases of falsidical paradoxes, such as the two-envelope
592  paradox. Here, there is a clear research result: RCT needs revision.

593  Third, there are some clear cases of apparent paradoxes, such as the self-torturer, in which
594  the whole bluster is caused by a fallaciously set up argument.

595  Finally, there are cases on which researchers cannot agree. These include Newcomb's
596  problem, Allais' and Ellsberg's paradox, which vacillates between veridical and falsidical
597  assessment, and the Backward-induction paradox, which vacillates between falsidical and
598  apparent assessment. In all these cases the verdict is still open as to whether they necessitate
599  RCT revision or not. Hence their continuing examination is part of active research in this area.

## References

Au4  600

602  Allais M (1953) Le comportement de l'homme rationnel
603  devant le risque: critique des postulats et axiomes de
604  l'école Américaine. Econometrica 21:503–546
605  Arntzenius F, McCarthy D (1997) Self torture and group
606  beneficence. Erkenntnis 47(1):129–144
607  Aumann R (1995) Backward induction and common
608  knowledge of rationality. Game Econ Behav 8:6–19
609  Bicchieri C (1989) Self refuting theories of strategic
610  interaction: a paradox of common knowledge.
611  Erkenntnis 30:69–85
612  Bicchieri C (1993) Rationality and coordination.
613  Cambridge University Press, Cambridge
614  Campbell R, Sowden L (eds) (1985) Paradoxes of ratio-
615  nality and cooperation: prisoner's dilemma and
616  Newcomb's problem. University of British Columbia
617  Press, Vancouver
618  Chew SH (1983) A generalization of the quasilinear mean
619  with application to the measurement of income
620  inequality and decision theory resolving the Allais
621  paradox. Econometrica 51:1065–1092
622  Clark M, Shackel N (2000) The two-envelope paradox.
623  Mind 109(435):415–442
624  Daniels N (2008) Reflective equilibrium. In: Zalta EN
625  (ed) The Stanford encyclopedia of philosophy,
626  fallth edn. The Metaphysics Research Lab, Stanford,
627  Available online: http://plato.stanford.edu/archives/
628  fall2008/entries/reflective-equilibrium/
629  Diekmann A, Mitter P (1986) Paradoxical effects of social
630  behavior: essays in honor of Anatol Rapoport.
631  Physica-Verlag, Heidelberg and Vienna, Available
632  online: http://www.socio.ethz.ch/vlib/pesb/index
633  Dietrich F, List C (2005) The two-envelope paradox: an
634  axiomatic approach. Mind 114:239–248
635  Ellsberg D (1961) Risk, ambiguity, and the savage axioms.
636  Q J Econ 75(4):643–669
637  Fox CR, Tversky A (1995) Ambiguity aversion and com-
638  parative ignorance. Q J Econ 110(3):585–603
639  Gauthier D (1986) Morals by agreement. Oxford Univer-
640  sity Press, Oxford

641  Gibbard A, Harper W (1981) Counterfactuals and
642  two kinds of expected utility. In: Harper W,
643  Stalnaker R, Pearce G (eds) Ifs: conditionals, belief,
644  decision, chance, and time. Reidel, Dordrecht,
645  pp 153–190
646  Guala F (2000) The logic of normative falsification: ratio-
647  nality and experiments in decision theory. J Econ
648  Methodol 7(1):59–93
649  Hansson SO, Grüne-Yanoff T (2009) Preferences. In:
650  Zalta EN (ed) The stanford encyclopedia of philos-
651  ophy, springth edn. The Metaphysics Research
652  Lab, Stanford, Available online: http://plato.stanford.
653  edu/entries/preferences/
654  Hardin G (1968) The tragedy of the commons. Science
655  162(3859):1243–1248
656  Hargreaves-Heap S, Hollis M, Lyons B, Sugden R, Weale
657  A (1992) The theory of choice: a critical guide.
658  Blackwell, Oxford
659  Horgan T (1985) Counterfactuals and Newcomb's prob-
660  lem. In: Campbell R, Sowden L (eds) Paradoxes of
661  rationality and cooperation: Prisoner's dilemma and
662  Newcomb's problem. University of British Columbia
663  Press, Vancouver, pp 159–182
664  Horwich P (1987) Asymmetries in time. MIT Press,
665  Cambridge, MA
666  Hyde D (2008) Sorites paradox. In: Zalta EN (ed) The
667  stanford encyclopedia of philosophy, fallth edn.
668  The Metaphysics Research Lab, Stanford, Available
669  online: http://plato.stanford.edu/archives/fall2008/
670  entries/sorites-paradox/
671  Jeffrey R (1990) The logic of decision, 2nd edn. University
672  of Chicago Press, Chicago
673  Joyce J (1999) The foundations of causal decision theory.
674  Cambridge University Press, Cambridge
675  Kahneman D, Tversky A (1979) Prospect theory: an
676  analysis of decision under risk. Econometrica
677  47:263–291
678  Koons R (1992) Paradoxes of belief and strategic ratio-
679  nality. Cambridge University Press, Cambridge

Kuhn S (2009) Prisoner's dilemma. In: Zalta EN (ed) The Stanford encyclopedia of philosophy, springth edn. The Metaphysics Research Lab, Stanford, Available online: http://plato.stanford.edu/archives/spr2009/entries/prisoner-dilemma/

Lewis D (1979) Prisoner's dilemma is a Newcomb problem. Philos Public Aff 8:235–240

Luce RD, Raiffa H (1957) Games and decisions: introduction and critical survey. Wiley, New York

Lycan WG (2010) What, exactly, is a paradox? Analysis 70(4):615–622

Mas-Colell A, Whinston MD, Green JR (1995) Microeconomic theory. Oxford University Press, New York

Meacham CJG, Weisberg J (2003) Clark and Shackel on the two-envelope paradox. Mind 112(448):685–689

Nozick R (1969) Newcomb's problem and two principles of choice. In: Rescher N (ed) Essays in honor of Carl G. Hempel. Reidel, Dordrecht, pp 114–146

Pettit P, Sugden R (1989) The backward induction paradox. J Philos 86:169–182

Quiggin J (1982) A theory of anticipated utility. J Econ Behav Organ 3(4):323–343

Quine WVO (1966) The ways of paradox. In: The ways of paradox and other essays. Random House, New York, pp 3–20, Original published as Paradox (1962) in Scientific American 206(4):84–95

Quinn WS (1990) The puzzle of the self-torturer. Philos Stud 59(1):79–90

Resnik MD (1987) Choices: an introduction to decision theory. University of Minnesota Press, Minneapolis

Richmond C, Sowden L (eds) (1985) Paradoxes of rationality and cooperation: Prisoners' dilemma and Newcomb's problem. University of British Columbia Press, Vancouver

Sainsbury RM (1988) Paradoxes. Cambridge University Press, Cambridge

Savage LJ (1954) The foundations of statistics. Wiley, New York

Schmeidler D (1989) Subjective probability and expected utility without additivity. Econometrica 57:571–587

Skyrms B (1980) Causal necessity: a pragmatic investigation of the necessity of laws. Yale University Press, New Haven

Vanderschraaf P, Sillari G (2009) Common knowledge. In: Zalta EN (ed) The stanford encyclopedia of philosophy, springth edn. The Metaphysics Research Lab, Stanford, Available online: http://plato.stanford.edu/archives/spr2009/entries/common-knowledge/

von Neumann J, Morgenstern O (1947) The theory of games and economic behavior, 2nd edn. Princeton University Press, Princeton

Voorhoeve A, Binmore K (2006) Transitivity, the Sorites paradox, and similarity-based decision-making. Erkenntnis 64(1):101–114

vos Savant M (1990) Ask marilyn column. Parade magazine, 9 Sept 1990, p 16

# Author Query Form

**Handbook of Risk Theory**
**Chapter No.: 27**

| Query Refs. | Details Required | Author's response |
|---|---|---|
| AU1 | Please check if insertion of closing quote is okay in extracted text. | |
| AU2 | Please confirm if the inserted in-text citations for Tables 27.3 and 27.4 are correct. | |
| AU3 | Please provide a caption for Figure 1. | |
| AU4 | Kindly cite the following references in text: Campbell and Sowden (eds) 1985; Lewis 1979. | |