

This is a section of [doi:10.7551/mitpress/11252.001.0001](https://doi.org/10.7551/mitpress/11252.001.0001)

The Handbook of Rationality

Edited by: Markus Knauff, Wolfgang Spohn

Citation:

The Handbook of Rationality

Edited by: Markus Knauff, Wolfgang Spohn

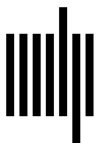
DOI: 10.7551/mitpress/11252.001.0001

ISBN (electronic): 9780262366175

Publisher: The MIT Press

Published: 2021

Funding for the open access edition was provided by the MIT Libraries Open Monograph Fund.



The MIT Press

8.2 Standard Decision Theory

Martin Peterson

Summary

What is it rational to do in a choice situation, given the limited and sometimes unreliable information available to an agent? The dominant view among decision theorists is that rational agents act in accordance with the principle of maximizing expected utility. This chapter surveys the debate on the expected utility principle and discusses the most significant objections raised in the literature.

1. Right versus Rational Decisions

Decision theory is a theory of practical rationality that seeks to answer the following question: what is it rational to do in a choice situation, given the limited and sometimes unreliable information available to an agent? Decision theorists distinguish between *right* decisions and *rational* ones. If you play the National Lottery and win, your decision was right (because the outcome was optimal), but if the odds of winning were as long as they usually are (and the act of playing the lottery did not in itself give you great pleasure), the decision was irrational. In light of the information available to you, you had more reason to abstain from playing than for doing so: the value of hitting the jackpot did not outweigh the low probability of winning. The task decision theorists set themselves is to articulate and defend this line of thought. How, exactly, should information about an act's possible outcomes guide our decisions? This chapter surveys some of the most influential answers and discusses a number of significant objections raised in the literature.

A fundamental assumption accepted by more or less all decision theorists is that the notion of rationality described by decision theory is a means–ends notion in which the rationality of a decision depends on the agent's beliefs and desires (see chapter 2.1 by Broome, this handbook). Decision theorists have little to say about what beliefs and desires are, but they assume that beliefs and desires can vary in strength. It is widely believed that the agent's degree of belief in something

can be represented by a subjective probability function (see section 4 of this chapter and chapter 4.1 by Hájek & Staffel, this handbook) and that desires can be represented by a utility function. The stronger you believe something and the more you desire it, the higher is your subjective probability and utility.

2. The Principle of Maximizing Expected Utility

Pascal (1670/1967), Nicolaus Bernoulli (in a letter from 1713), and Cramer (1728/1975) propose that an act is rational if and only if its *expected utility* is at least as high as that of every alternative act.¹ This has been the dominant view among decision theorists ever since, but there is significant disagreement about how the expected utility principle should be interpreted and justified. The expected utility of an act with n possible outcomes x_1, \dots, x_n is

$$\sum_{i=1}^n p_i \cdot u(x_i), \quad (1)$$

where p_i is the probability of outcome x_i and $u(x_i)$ the utility of x_i (see chapter 8.1 by Grüne-Yanoff, this handbook). If the number of possible outcomes x is infinite, the expected utility is $\int_{x \in X} [p(x) \cdot u(x)] dx$.

3. The Law of Large Numbers

A possible reason for accepting the principle of maximizing expected utility (or expected value more generally) is that agents who do so will almost certainly be better off in the long run. This insight is formalized in a mathematical theorem known as the law of large numbers:² if a random experiment for which the expected value is E_x is repeated n times, then the probability that the average value \bar{X} of the experiment differs from E_x by more than some small number ε converges to 0 as the number of experiments n goes to infinity, and this holds true for every positive $\varepsilon > 0$ no matter how small:

$$\lim_{n \rightarrow \infty} p(|\bar{X} - E_x| > \varepsilon) = 0. \quad (2)$$

This mathematical insight might be of some importance for agents who face similar choice situations a large number of times. However, many decisions are unique. Claims about what *would* happen if one were to make decisions that one will *not actually* make therefore seem irrelevant. Would it, for instance, be rational for you (given your interests and financial circumstances) to study at university *X* rather than university *Y*? Should you buy the house offered for sale for \$400,000, or should you look for a cheaper house? Would it be rational to quit your job to sail around the world? When answering these questions, it seems irrelevant what would happen if you were to face similar choices a million times, because you *know* that you will only face each choice a small number of times. Keynes (1923/2000) expressed this insight in its most brutal form: “In the long run we are all dead” (p. 80).

4. Axiomatic Analyses of the Expected Utility Principle

4.1 Von Neumann and Morgenstern

Several authors have developed axiomatic analyses of the expected utility principle, which are sometimes interpreted as normative arguments for accepting the principle. See, for example, Ramsey (1926/1931), von Neumann and Morgenstern (1947), Savage (1954/1972), and Jeffrey (1983).³ None of these axiomatizations entail that rational agents will make their decisions by consciously multiplying probabilities and utilities. The axioms merely entail that agents forming preferences *in accordance with* the axioms can be described *as if* they were maximizing expected utility. As Blackburn (1998) puts it, these representation and uniqueness theorems provide us with a “grid for imposing interpretation: a mathematical structure, designed to render processes of deliberation mathematically tractable, whatever those processes are” (p. 135).

Von Neumann and Morgenstern focus on decision problems in which probabilities can be defined exogenously, meaning that the agent somehow knows the probability of every event relevant to her decision. In this section, we will use a slightly different notation than that used in chapter 8.1 by Grüne-Yanoff (this handbook). Recall that von Neumann and Morgenstern conceive of acts as lotteries: *xpy* is the lottery that yields *x* with probability *p* and *y* with probability $1-p$ (and this holds regardless of whether *x* and *y* are basic prizes or new lotteries). The possible outcomes *x* or *y* are either basic prizes (such as money or a ticket to the opera) or lotteries over basic prizes. The relation \succsim is a binary preference relation over the set of all possible lotteries *L*. Von Neumann and Morgenstern propose that agents are

rationally required to form preferences over lotteries that meet the following structural requirements, for all lotteries and basic prizes *x, y, z* and all probabilities *p* and *q*:

- vNM 1 (Weak Order) \succsim is a complete, reflexive, and transitive relation on *L*.
- vNM 2 (Independence) $x \succsim y$ if and only if $xpz \succsim ypz$.
- vNM 3 (Continuity) If $x \succ y \succ z$, then there exist some *p* and *q* such that $xpz \succ y \succ xqz$.

The independence axiom is the most controversial of the three axioms. In table 8.2.1, it seems reasonable to prefer lottery *x* to *y*, because this will give you \$1 million for sure. However, the independence axiom entails that if you prefer *x* to *y*, then you must also prefer *xpz* to *ypz*. This is counterintuitive. It seems perfectly reasonable to prefer *ypz* to *xpz*, because there is no certain outcome here. The .10 probability of winning \$5 million instead of \$0 might very well be considered more attractive than a slightly larger probability (.11) of winning a much smaller amount (\$1 million). For discussions of the ordering and continuity axioms, see chapter 8.1 by Grüne-Yanoff (this handbook).

Von Neumann and Morgenstern’s theorem states that whenever vNM 1–3 are satisfied, there exists a real-valued function *u(x)* that takes a lottery in *L* as its argument and returns a real number between 0 and 1, which has the following properties:

- (1) $x \succ y$ if and only if $u(x) > u(y)$.
- (2) $u(xpy) = p \cdot u(x) + (1 - p) \cdot u(y)$.
- (3) For every other function *u'* satisfying (1) and (2), there exist some constants *c* > 0 and *d* such that $u' = cu + d$.

Properties (1) and (2) constitute the representation part: the agent assigns higher utility to outcomes she judges to be preferable, and the utility of a lottery equals the expected utility of its outcomes. Property (3) is the uniqueness part: any other function *u'* that satisfies the first two properties

Table 8.2.1
The independence axiom

	Ticket 1	Ticket 2–11	
Lottery <i>x</i>	\$1 million	\$1 million	
Lottery <i>y</i>	\$0	\$5 million	
	Ticket 1	Ticket 2–11	Ticket 12–100
Lottery <i>xpz</i>	\$1 million	\$1 million	\$0
Lottery <i>ypz</i>	\$0	\$5 million	\$0

is a linear transformation of u , meaning that utility is measured on an interval scale (see chapter 8.1).

A common objection to normative interpretations of von Neumann and Morgenstern’s result is that it puts the cart before the horse (for an overview, see Peterson, 2008, section 2.5). The agent acting in accordance with these axioms does not prefer an act *because* its expected utility is optimal. She can be described *as if* she were acting from the expected utility principle, but we have no reason to believe that the utility function we ascribe to the agent is equivalent to the function that guides her choices or that such a function exists. The problem is thus that the theory is not action-guiding. The ordering axiom explicitly requires that agents have complete preferences over *all possible* lotteries (acts), including very complex lotteries (acts), such as the ones mentioned earlier: Would it be rational for you to study at university X rather than university Y ? Should you buy the house offered for sale for \$400,000, or should you wait and hope the price drops to \$350,000? Should you quit your job to sail around the world? Agents who have complete preferences over all possible lotteries hardly need advice from decision theorists.

The standard response to this objection is that the axioms could be action-guiding in an indirect sense. They can be treated as “inference rules” that can help nonideal agents revise rationally impermissible preferences. If you, say, discover that your preferences are incomplete, or intransitive, then the axioms tell you how to overcome this deficiency: revise your preferences, in whatever way you wish, so they meet the structural requirements imposed by the axioms. However, this type of indirect action guidance seems somewhat unhelpful. If you are trying to figure out if it is rational for you to study at university X or university Y or whether to offer \$400,000 for a house you like, and the decision theorist is merely able to tell you that “anything goes” *as long as the axioms are fulfilled*, then decision theory will not be of much help to you.

4.2 Savage

The axiomatizations proposed by Ramsey (1926/1931), Savage (1954/1972), and Jeffrey (1983) differ from von Neumann and Morgenstern’s in that probabilities are derived from the axioms themselves. In these axiomatizations, the probability function is not an exogenously defined function. Furthermore, since the axioms regulate the agent’s preferences, we can interpret probabilities derived from these axioms as *subjective* probabilities, or credences (see chapter 4.1 by Hájek & Staffel, this handbook). To say that the probability is .90 that it will rain tomorrow means that you believe to a fairly strong

degree that it will rain tomorrow: your credence makes you willing to pay something worth .90 units of utility for a bet in which you win 1 unit if it rains tomorrow and nothing otherwise.

Savage conceives of acts as functions from a set of possible states of the world to outcomes. The act of not bringing the umbrella to the football game is a function in which the two possible states are “rain” and “no rain,” and the possible outcomes are “the agent gets wet” and “the agent does not get wet,” depending on which state turns out to be the true state of the world.

Some of the axioms proposed by Savage resemble those proposed by von Neumann and Morgenstern. This includes the assumption that the agent’s preferences satisfy the weak ordering axiom, which raises the same issue about putting the cart before the horse discussed earlier. Another axiom borrowed from von Neumann–Morgenstern (and, ultimately, from Ramsey, 1926/1931) is the so-called *sure-thing principle*. This is Savage’s name for a much-discussed version of the independence axiom. In Savage’s terminology, two acts f and g agree with each other in a set S of states if and only if $f(s) = g(s)$ for all $s \in S$. In this terminology, the sure-thing principle can be formulated as follows:

If f and g , and f' and g' , are such that

1. in $\neg S$, the complement of S , f agrees with g , and f' agrees with g' ,
2. in S , f agrees with f' , and g agrees with g' , and
3. $f \succ g$,

then $f' \succ g'$.

Table 8.2.2 is helpful for highlighting the similarities with von Neumann and Morgenstern’s independence axiom. Imagine that you are offered a choice among four lotteries, each of which has 100 tickets. It seems rationally permissible to prefer f to g but g' to f' (for the reason mentioned in the discussion of the independence axiom), but the sure-thing principle controversially entails that if you prefer f to g , then you must prefer f' to g' .

Table 8.2.2

The sure-thing principle

	S		$\neg S$
	Ticket no. 1	Ticket nos. 2–11	Ticket nos. 12–100
Act f	\$1 million	\$1 million	\$1 million
Act g	\$0	\$5 million	\$1 million
Act f'	\$1 million	\$1 million	\$0
Act g'	\$0	\$5 million	\$0

One of the axioms in Savage's theory, which does not figure in von Neumann and Morgenstern's (in which the probability function is defined exogenously), is designed to test whether the agent considers two states to be equally probable. Savage's own formulation of this axiom is quite complex, but the following explanation offered by Broome (1990) is more helpful:

Savage . . . needs to test whether two events, say E and F , are equally probable. He does this by taking a pair of outcomes, say A and B , that are known not to be indifferent. He forms a gamble $(A, E; B, F)$ in which A comes about in event E and B in F . And he forms the opposite gamble $(B, E; A, F)$. The events are equally probable if and only if these gambles are indifferent. Savage's Postulate 4 says that this test will deliver the same answer whatever pair of non-indifferent outcomes A and B are used. (p. 483)

Broome notes that a problem with this axiom is that it seems to assume what it is supposed to show: the assumption that the agent will be indifferent between the gambles only if they are judged to be equally probable *presupposes* that probabilities and utilities are aggregated in accordance with the expected utility principle.

Savage's uniqueness and representation theorem is similar to von Neumann and Morgenstern's, except that the probability function is derived from the axioms. The theorem guarantees the existence of a utility function that is unique up to positive linear transformations and a probability function such that one act is preferred to another just in case its expected utility exceeds that of the latter.

4.3 Bolker–Jeffrey

The theory proposed in Jeffrey (1965) was axiomatized by Bolker (1967) and presented in a slightly modified form in the second edition of Jeffrey's book (Jeffrey, 1983). The Bolker–Jeffrey theory differs from von Neumann and Morgenstern's and Savage's in several important respects.

In the Bolker–Jeffrey theory, it is assumed that the agent has preferences over *propositions*, such as "I have a refreshing swim with friends," rather than lotteries (von Neumann & Morgenstern, 1947), or functions from states to outcomes (Savage, 1954/1972). Utilities and probabilities are also assigned to propositions, meaning that the Bolker–Jeffrey theory is more economical from an ontological point of view. Another philosophical advantage is that the use of propositions enables the agent to retain her original beliefs about the world throughout the deliberative process. There is no need for the agent to imagine complex lotteries that require the existence of causal processes she does not believe in. The agent can continue to believe whatever she believes, without adding any new beliefs or modifying old ones.

In the Bolker–Jeffrey theory, it is assumed that preferences over the Boolean algebra B (with the zero removed; contradictions have no place in the preference ordering) satisfy the weak ordering axiom, are continuous, and satisfy the following two axioms:

- (Averaging) If a, b in B are contraries (i.e., cannot both be true), then
 $a \succ b$ implies $a \succ (a \vee b) \succ b$, and
 $a \sim b$ implies $a \sim (a \vee b) \sim b$.
- (Impartiality) If a, b, c in B are pairwise contraries, and $a \sim b$ but not $a \sim c$, and $(a \vee c) \sim (b \vee c)$, then for every d in B that is contrary to a and b , $(a \vee d) \sim (b \vee d)$.

The averaging axiom is somewhat similar to, but weaker than, the independence axiom. It requires that the agent ranks a disjunction of two propositions somewhere between the two disjuncts. This may seem reasonable. However, a possible reason for ranking a disjunction above or below the disjuncts could be that the uncertainty introduced by the disjunction sometimes has a positive or negative effect on the value of the proposition. Imagine, for instance, that you prefer the proposition "My partner loves me" to "My partner breaks up with me tomorrow." If so, it does not seem irrational to rank the disjunction of these two propositions lower than "My partner breaks up with me tomorrow." Not knowing for sure if your relationship is over might be worse than knowing that it is. If you learn that your partner will break up with you tomorrow, you can adjust your plans for the future and move on.

The impartiality axiom resembles the equiprobability axiom in Savage's theory. The intuition it seeks to capture is similar to that summarized by Broome in the quote in section 4.2: to test whether the agent considers two propositions to be equally probable, we can ask whether the agent is indifferent between $a \vee c$ and $b \vee c$ whenever she is indifferent between a and b . If the agent, for instance, is indifferent between a and b but prefers $a \vee c$ over $b \vee c$, then she does not consider a and b to be equally probable (as long as she maximizes expected utility; more about this soon). The impartiality axiom states that this holds true for every third proposition c .

Broome (1990) points out that the impartiality axiom assumes what it is supposed to show, just as axiom 4 in Savage's theory. If the agent prefers a to b , then in case she maximizes expected utility, the fact that she prefers $a \vee c$ over $b \vee c$ entails that she considers a to be more probable than b . However, if she aggregates probabilities and utilities in some other way, she has no reason to

accept this axiom. It is not entirely satisfactory to justify the expected utility principle by appealing to an axiom that will be attractive primarily to those who accept the expected utility principle.

4.4 Other Axiomatizations

A radically different approach to the expected utility principle is to start with preferences over certain outcomes (basic prizes), as well as beliefs about the world, and then generate preferences over risky acts from this *ex ante* information. The *ex ante* approach steers clear of the objection that traditional axiomatizations put the cart before the horse. The drawback is that it seems somewhat optimistic to think that ordinary agents have access to exogenously defined utility and probability functions (see Peterson, 2004).

5. Pragmatic Arguments

Faced with counterexamples to some of the key axioms of expected utility theory, decision theorists have attempted to offer further warrant for the axioms proposed by von Neumann and Morgenstern, Savage, and others by formulating *pragmatic arguments* (see, e.g., McClennen, 1990). A pragmatic argument is a device that seeks to show that anyone who violates a purported rationality requirement can be placed in a situation in which she has to act contrary to her own preference.

Donald Davidson, J. C. C. McKinsey, and Patrick Suppes (1955) propose the following pragmatic argument in defense of the transitivity condition, which is key to the assumption that rational preferences form a weak order:

Mr. S. is offered his choice of three jobs by a cynical department head (never mind what department): He can be a full professor with a salary of \$5,000 (alternative *a*), an associate professor at \$5,500 (alternative *b*) or an assistant professor at \$6,000 (alternative *c*). Mr. S. reasons as follows: $a \succ b$ since the advantage in kudos outweighs the small difference in salary; $b \succ c$ for the same reason; $c \succ a$ since the difference in salary is now enough to outweigh a matter of rank. . . . It is clear that the set of Mr. S.'s preferences makes a rational choice impossible, for whichever alternative he chooses there will be another alternative which is preferred to it. (Davidson et al., 1955, p. 145)

This argument is sometimes presented in a *diachronic* form, that is, as a series of choices to be made at different points in time. Imagine, for instance, that you sign the contract for the assistant professorship with a salary of \$6,000 (alternative *c*). Right after you do so, the cynical department head reminds you that you prefer *b* to *c*. He offers you to swap if you pay him a small administrative

fee, say \$.01. Because you prefer *b* minus the small fee to *c*, you accept his offer. However, right after you sign the contract for job *b*, he offers you to swap *b* for *a* if you pay a small fee, and once you have paid the fee, he offers you to pay him for swapping again, and so on. The cyclical structure of your preferences forces you to keep paying small amounts for swapping over and over again. After 100 million swaps, you have lost \$1 million but gained nothing, which seems to indicate that your preferences are irrational. This type of pragmatic argument is known as a *money-pump argument*.

Schick (1986) claims that a rational agent with cyclical preferences will “see which way the wind is blowing . . . [and] seeing what is in store for him, he may well reject the offer and thus stop the pump” (p. 117). This is, however, a problematic response. Schick's argument can at best cast doubt on the diachronic money pump sketched above, but it does not affect the original *synchronic* version proposed by Davidson et al. (1955). As stressed by Gustafsson (2013, p. 462), Davidson et al.'s point is that “it is irrational to choose an alternative to which another alternative is preferred.” This is irrational even if the agent makes only *one* such choice. There is no need to consider diachronic versions of the argument in which the agent has to take into account what will happen in the future.

Another problem with money-pump arguments for transitivity is that they do *not* show that *every possible* violation of transitivity will lead to a sure loss. Preferences can violate the transitivity condition in many different ways, and not every violation forces the agent to act against her own preference. Imagine, for instance, that you strictly prefer *a* to *b*, *b* to *c*, but are indifferent between *a* and *c*. With these preferences, you would be rationally obliged to pay a small amount for swapping *c* for *b* and then pay again for swapping *b* for *a*. However, once you have acquired *a*, you have no reason to continue swapping. Since you are indifferent between *a* and *c*, you may keep *a*, which stops the pump. By keeping *a*, you do not choose an alternative to which another is preferred. However, because you are indifferent between *a* and *c*, your preference violates the transitivity condition.

The worry outlined above can be further analyzed by distinguishing between strong (“forcing”) and weak (“nonforcing”) pragmatic arguments (see, e.g., Peterson, 2015). In a *strong* pragmatic argument, the agent's preferences *guarantee* that a rational agent will choose an alternative to which another is preferred. In a *weak* pragmatic argument, the agent's preferences *permit* the rational agent to choose an alternative to which another is preferred. If you, for instance, strictly prefer *a* to *b*, and *b* to *c*, but are indifferent between *a* and *c*,

then you are rationally *permitted* to swap over and over again, each time paying a small fee when you swap to a strictly preferred option. However, you would not be rationally *obliged* to keep swapping once you reach *a*. Some scholars think that weak pumps are strong enough: by merely showing that your preferences *permit* you to be money-pumped, we can conclude that your preferences are irrational. Defenders of strong pumps insist that for a pragmatic argument to be successful, we have to establish that the agent is rationally *obliged* to be money-pumped.

The distinction between weak and strong pumps is also relevant for assessing pragmatic arguments for the assumption that rational preferences are complete. It is reasonable to think that if you have no preference between *a* and *c*, then you are *rationally permitted*, but not *rationally obliged*, to swap. Hence, if you strictly prefer *a* to *b*, and *b* to *c*, you would be vulnerable to a weak, but not a strong, money pump.

Some authors have argued that we can justify the independence axiom (see section 4.1) by appealing to pragmatic considerations. Figure 8.2.1 summarizes a sequential choice problem with two choice nodes (boxes) and three chance nodes (circles). At the first choice node (the leftmost box), the agent chooses between going “up” and “down.” If the agent goes up and event $\neg E$ occurs (the probability of which is $89/100$), she wins \$0. However, if event E occurs (the probability for this is $11/100$), the agent faces a new choice, in which she can choose to either walk away with \$1 million or accept a lottery in which she wins

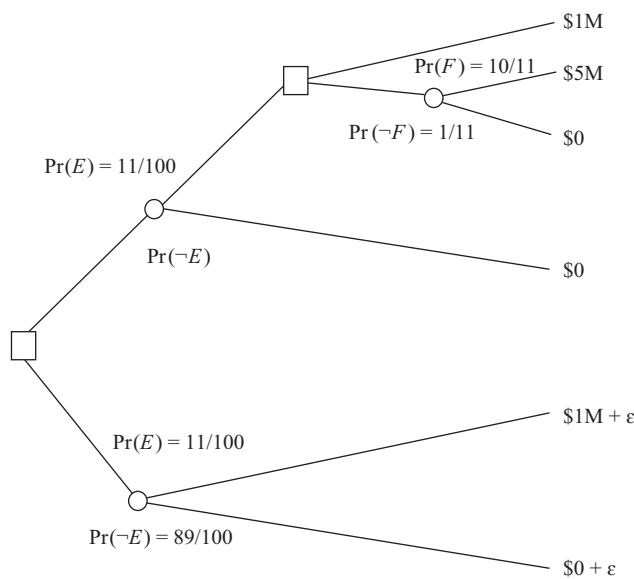


Figure 8.2.1

A decision tree (after Rabinowicz, 1995).

\$5 million with probability $10/11$ and \$0 with probability $1/11$. If the agent goes down at the first choice node, she faces a lottery in which the outcomes and probabilities are exactly as if she were to decide at the outset to go up at both the first and second choice nodes, except that a small bonus ε has been added to each possible outcome.

The outcomes and probabilities in figure 8.2.1 are analogous to those in table 8.2.1. Note that

x : \$1 million;

y : \$5 million with probability $10/11$, otherwise \$0;

xpz : \$1 million with probability $11/100$, otherwise \$0;

ypz : \$5 million with probability $10/100$, otherwise \$0.

To construct a pragmatic argument for the independence axiom, suppose the agent's preferences are as follows:

$$x \succ y, \quad (1)$$

$$ypz \succ xpz. \quad (2)$$

For the reason explained in section 4.1, these preferences violate the independence axiom. To construct the argument, we also assume that the small bonus ε in figure 8.2.1 does not reverse the agent's preference between ypz and xpz . Given that the bonus is small enough, this is not a controversial assumption:

$$ypz \succ xpz + \varepsilon \succ xpz. \quad (3)$$

Because the agent prefers ypz to xpz , per (2), she will go up at the first choice node in figure 8.2.1 with the intention to go down at the second if event E occurs. However, if the agent reaches the second choice node, she will go up, since she prefers x to y , per (1). According to the agent's preferences, \$1 million for sure is worth more to her than the risky gamble in which she wins \$5 million with probability $10/11$. However, this means that she is willing to forsake a small but certain bonus. If the agent goes down at the first choice node, she would obtain the same outcome *plus* the bonus ε . From a pragmatic point of view, agents who act in accordance with (1)–(3) thus forsake sure gains, compared to what they could have achieved had they instead revised their preferences in accordance with the independence axiom.

Is this a convincing pragmatic argument for the independence axiom? To start with, note that the argument only purports to show that *one of many possible violations* of the independence axiom causes the agent to forsake the small bonus ε . The argument does not purport to show that *every possible violation* has the same pragmatic effect. Suppose, for instance, that the agent prefers x to y but is indifferent, or has no preference, between ypz and xpz . If this is the case, she is at most vulnerable to a weak pragmatic argument but not to a strong one.

A more fundamental concern is that the argument presupposes a theory of sequential decision-making few people find plausible. In essence, the problem is that the agent who goes up at the first choice node with the intention of going down at the second is overly *myopic* (shortsighted). Schick's (1986, p. 117) observation that a rational agent will see "what is in store for him" seems equally applicable here as in the money-pump argument for transitivity mentioned earlier. According to Schick, a *sophisticated* agent facing a sequential choice problem should reason backward: if the agent were to reach the second choice node, she would prefer x to y , and she knows this already at the first choice node. Therefore, there are only two options open to her: she can (i) go down at the first choice node or (ii) go up at the first choice node with the intention to go up at the second. The third alternative, up followed by down, is not a genuine alternative. It should thus be deleted from the list of options. Therefore, the example in figure 8.2.1 fails to show that agents who violate the independence axiom will act contrary to their own preference.

Rabinowicz (1995) argues that it is possible to overcome this objection: the sophisticated approach does not make the agent immune to pragmatic arguments. In figure 8.2.2, the sophisticated agent with the preferences described in (1) and (2) would reason backward just as before: she prefers \$1 million for sure (i.e., x) to \$5 million with a probability of 10/11 (y), so she knows already at the first choice node that she would go up at the second choice node. At the first node, she

thus faces a choice between (i) going down and (ii) going up with the intention of going up at the second choice node. Because the agent prefers ypz to xpz , and the subtraction of ϵ from the outcomes of ypz does not alter her preference, she will go down at the first node. However, according to Rabinowicz, down is dominated by the plan to go up at the first node and down at the second. From a pragmatic point of view, the sophisticated agent is thus acting contrary to her own preference. She prefers ypz to $ypz - \epsilon$, but the sophisticated approach forces her to choose $ypz - \epsilon$. Advocates of the sophisticated approach respond to this that the dominating plan (up, down) is not feasible and can therefore be ignored.

McClennen (1990) rejects the myopic and sophisticated approaches and proposes what he calls a *resolute* approach. Resolute agents adopt a plan at the first choice node for how to act at all later choice nodes and then stick to the initial plan as long as no new information is received. The plan will typically be conditional: "Go up at the second node if event E occurs" and so on. Resolute agents who violate the independence axiom cannot be money-pumped, but a common objection is that resolute choice presupposes perfect self-control (for a discussion, see, e.g., Peterson & Vallentyne, 2018).

6. Causal versus Evidential Decision Theory

Decision theorists disagree on how rational agents should respond to information about causal processes. The *locus classicus* for this debate is a thought experiment first discussed in print by Nozick (1969), who attributes it to the physicist William Newcomb.

Imagine a supercomputer that is very good at predicting people's choices by scanning their brains. You are asked to choose between two boxes. The first box is transparent, and you can see that it contains exactly \$1,000. The other box is not transparent. It contains either \$0 or \$1 million, depending on how the supercomputer predicts you will choose. The operator, whom you have no reason to distrust, tells you that (i) if the supercomputer predicts that you will take only the first box, then the operator will put \$1 million in the second box, but (ii) if the supercomputer predicts that you will take both boxes, then the operator will put \$0 in the second box. The supercomputer's prediction is highly accurate but not infallible. Of the 1,000 predictions the computer has made so far, the computer made correct predictions 990 times.

You are now asked to choose between the following options:

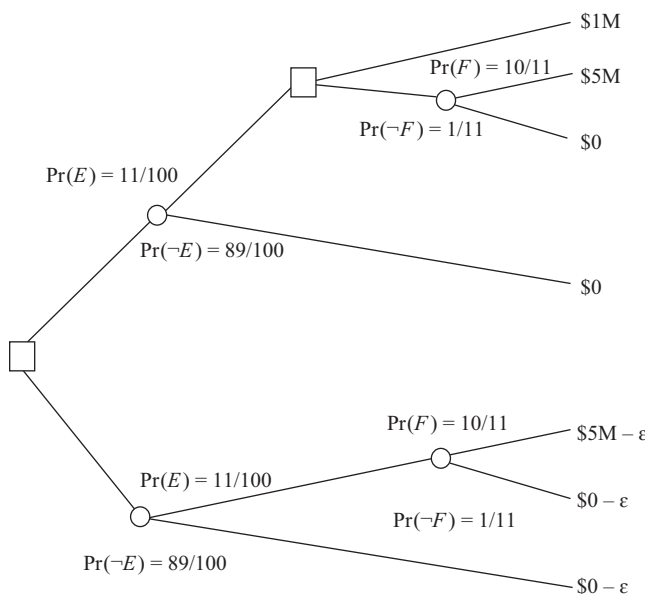


Figure 8.2.2 A decision tree for sophisticated agents (after Rabinowicz, 1995).

- (a) Take box 1 (\$1,000) and box 2 (either \$0 or \$1 million).
 (b) Take only box 2 (either \$0 or \$1 million).

Note that the supercomputer *first* predicts how you will choose; *then* the operator adjusts the content in box 2 by putting either \$1 million or \$0 in the box; and *finally*, after these events have already happened, you make your choice. The temporal order is important because it guarantees that your choice *does not causally influence* the content of box 2.

Causal decision theory (CDT) holds that the only rational option is to take both boxes. Because the operator has already put \$1 million or \$0 in the second box when it is time for you to choose, the decision you make will not affect the content of the second box. If you take both boxes, you will thus get \$1,000 in addition to whatever is in the second box. Another way to put this is to say that taking two boxes *dominates* taking one box: no matter what the world is like (\$0 or \$1 million in the second box), you will be at least as well or better off if you take both boxes (for a detailed defense of CDT, see, e.g., Joyce, 1999).

Advocates of evidential decision theory (EDT) take issue with what they consider to be an overly narrow analysis offered by causal decision theorists. Their point is that the analysis offered by causal decision theorist does not accurately reflect the relevant information available to the agent. In Newcomb's problem, you know that the probability that the supercomputer's prediction is correct is 99%. This holds true even if you take only one box. The choice you make thus *counts as evidence* for what the supercomputer has predicted you will do. The only rational option according to advocates of EDT is to calculate the expected utility such that probabilities are conditional on which act is performed (see table 8.2.3).

To simplify the calculation, we assume that the agent's utility for money is linear. (Note that we could easily modify the prizes to fit any utility function the agent may have.) It can then be easily verified that the expected utility of taking only the second box exceeds that of taking both boxes, meaning that this is the only rationally permissible option.

Which theory do we have most reason to accept? Egan (2007) argues that the following example shows that CDT sometimes has unacceptable implications:

Paul is debating whether to press the "kill all psychopaths" button. It would, he thinks, be much better to live in a world with no psychopaths. Unfortunately, Paul is quite confident that only a psychopath would press such a button. Paul very strongly prefers living in a world *with* psychopaths to dying. Should Paul press the button? (p. 97)

Joyce (2012), an influential advocate of CDT, argues that "CDT does not recommend shooting as the uniquely rational act. . . . Agents who think their way through [this example] from the perspective of CDT will wind up being correctly indifferent between shooting and refraining once they have taken *all their causally relevant information in to account*" (p. 125). It is an open question whether this is a satisfactory response to Egan-style examples. Many people seem convinced that it is not rationally permissible to shoot or press the button (see also Ahmed, 2014, and chapter 7.1 by Pearl, this handbook).

7. Paradoxes and Puzzles

This section reviews some paradoxes and puzzles not mentioned in the previous sections.

7.1 Allais's Paradox

The point of departure for Allais's (1953) paradox, which is not a paradox in a strict philosophical sense but rather an unintuitive implication of the expected utility principle, is the four acts f , g , f' , and g' listed in table 8.2.2 in section 4.2. To state the paradox, we calculate the *arithmetic difference* in expected utility between acts f and g , as well as between f' and g' . Since the agent's utility of money is unknown, we write

$$\begin{aligned} u(f) - u(g) &= \\ &= u(1M) - [.01 \cdot u(0M) + .1 \cdot u(5M) + .89 \cdot u(1M)] \\ &= .11 \cdot u(1M) - [.01 \cdot u(0) + .1 \cdot u(5M)], \\ u(f') - u(g') &= \\ &= [.11 \cdot u(1M) + .89 \cdot u(0)] - [.9 \cdot u(0M) + .1 \cdot u(5M)] \\ &= .11 \cdot u(1M) - [.01 \cdot u(0) + .1 \cdot u(5M)]. \end{aligned}$$

As noted in section 4.2, it seems reasonable to prefer f to g but g' to f' , even though this violates Savage's sure-thing principle. However, the difference in expected utility between f and g is exactly the same as that between f'

Table 8.2.3
Newcomb's problem

	\$1 million in second box		\$0 in second box	
Take second box only	\$1 million	(prob. .99)	\$0	(prob. .01)
Take both boxes	\$1 million + \$1,000	(prob. .01)	\$1,000	(prob. .99)

and g' regardless of what the agent's utility of money is. Therefore, whenever f is preferred to g , the agent will also prefer f' to g' , regardless of the agent's utility of money. This challenge to the sure-thing principle is thus a challenge to the principle of maximizing expected utility. Regardless of how much or little you care about money, you are bound to violate the expected utility principle if you prefer f to g but g' to f' . The "paradox" consists in the fact that there is no way of assigning utility to money that renders the expected utility principle compatible with a preference for f over g and a preference for g' over f' .

The standard response to Allais's paradox is to point out that the description of the outcomes in table 8.2.2 is incomplete. It does not take into account that you will be *disappointed* if you choose act g and do not become a millionaire. A more accurate representation would be the one in table 8.2.4.

Although we do not know exactly what the agent's utility of "\$0 and disappointment" is, it seems reasonable to assume that it is less than the utility of \$0. Therefore, the difference in expected utility between f and g is not the same as that between f' and g' .

7.2 The St. Petersburg Paradox

The St. Petersburg paradox is one of the oldest and most well-known problems in decision theory. It is not a paradox in a strict philosophical sense (just like the other problems discussed here) but merely a counterintuitive implication of the expected utility principle.

Imagine that a fair coin is tossed n times until it lands heads up for the first time, after which you receive a prize worth 2^n units on your personal utility scale. If, for instance, the coin lands heads on the first toss, you win 2 units, and if you get to toss it five times, you win $2^5 = 32$ units, and so on. How much would a rational agent be willing to pay for an opportunity to playing this game?

The expected utility of the St. Petersburg game is

$$\frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 4 + \frac{1}{8} \cdot 8 + \dots = 1 + 1 + 1 + \dots = \infty.$$

Table 8.2.4
Allais's paradox

	Ticket 1	Ticket 2–11	Ticket 12–100
Act f	\$1 million	\$1 million	\$1 million
Act g	\$0 and disappointment	\$5 million	\$1 million
Act f'	\$1 million	\$1 million	\$0
Act g'	\$0	\$5 million	\$0

Hence, a rational agent would, according to the expected utility principle, be willing to sacrifice *all* her assets for the opportunity to play this game once. But this seems absurd. The most likely outcome of the game is that the agent wins 2, 4, or 8 units of utility. Can a remote possibility of winning a huge amount really compensate for a very likely loss of all one's assets?

In Bernoulli's (1738/1954) original formulation of the St. Petersburg paradox, the prizes were gold coins rather than units of utility. The concept of utility was invented by Cramer (1728/1975) as a way of explaining why a rational agent is not required to sacrifice all her assets for the opportunity to play the game. If the agent's marginal value (utility) of gold is decreasing, the expected utility of the gamble is finite. However, although the notion of decreasing marginal utility is still important in many economic analyses, modern philosophers are aware that the paradox can be reformulated in the manner mentioned above, in which the agent wins 2^n units of *utility* rather than some units of gold. This reformulation of the paradox of course presupposes that utility is unbounded (Peterson, 2020).

7.3 The Pasadena Puzzle

The gist of the Pasadena Puzzle is that some gambles seem to have no expected utility at all. Before we state the paradox, it is helpful to recall that some infinite series have no unique sum. Consider, for instance, the alternating harmonic series $\sum_n \frac{(-1)^{n-1}}{n}$. We can write this as $1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \frac{1}{5} \pm \dots$. The sum of this series, added up in this order, is $\ln 2$. However, the Riemann rearrangement theorem entails that the sum of the alternating harmonic series (and other conditionally convergent series) will depend on the *order* in which the terms are summed up. In fact, the terms can be rearranged in a permutation such that the alternating harmonic series converges to any given value, including $+\infty$ and $-\infty$.

To turn this somewhat esoteric mathematical fact into an objection to the expected utility principle, Nover and Hájek (2004) ask us to imagine that you toss a fair coin n times until it lands heads up for the first time. If n is odd, you win $2^n/n$ units of value, but if n is even, you have to *pay* $2^n/n$. Do we have any reason to think that the expected value of the Pasadena gamble is $\sum_n \frac{(-1)^{n-1}}{n} = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \frac{1}{5} \pm \dots = \ln 2$? Recall that we could also sum up the terms in some other order that would make the series converge to any real number we want. Suppose, for instance, that we write down each term of the sum on separate pieces of paper and drop all of them on the floor before we sum up the terms. Why

not just sum up the terms in the order we happen to pick them up from the floor?

Easwaran (2008) responds to this challenge by introducing a distinction between weak and strong expected value. He shows that the Pasadena gamble's (only) weak expected value is $\ln 2$. However, Lauwers and Vallyntyne (2016) remind us that not all gambles have any weak expected value. Their suggestion is that we should rather think of the entire interval of all possible values the series can take as its value. However, a problem with this proposal is that it is not clear how we ought to choose between something worth, say, $[-\infty, +\infty]$ and 0. The interval $[-\infty, +\infty]$ is clearly not worth as much as 0, because if it were, we could obviously reduce the value of the series to a point value, which would make the interval value redundant.

8. Conclusion

There is persistent disagreement about how the principle of maximizing expected utility should be interpreted and justified. None of the positions reviewed in this chapter appears to be fully satisfactory. However, despite all these problems, it seems likely that the expected utility principle will continue to dominate debates in decision theory for the years to come. There are few, if any, promising alternatives to the expected utility principle. Buchak's (2013) recent modification of the principle is worth considering, but the decision rule she proposes makes the agent vulnerable to pragmatic objections similar to those mentioned in section 5, as she points out herself. Kahneman and Tversky's (1979) prospect theory offers reasonable descriptive predictions, but no one thinks their theory would fare better than the expected utility principle from a normative perspective (see also chapter 8.4 by Hill, this handbook).

Notes

1. Pascal and Bernoulli discuss the principle of maximizing expected *monetary value*. Cramer's contribution was to clearly state the distinction between monetary value and utility. (He used the term "moral value.")
2. The law of large numbers comes in two versions, the "weak" and the "strong" law of large numbers. The version presented here is the weak version.
3. Note that the first edition of von Neumann and Morgenstern's book, published in 1944, does not include their axiomatization of the expected utility principle. The axioms first appeared in an appendix to the second edition published in 1947.

References

- Ahmed, A. (2014). *Evidence, decision and causality*. Cambridge, England: Cambridge University Press.
- Allais, M. (1953). Le comportement de l'homme rationnel devant le risque: Critique des postulats et axiomes de l'école Américaine [Rational man's behavior in the presence of risk: Critique of the postulates and axioms of the American school]. *Econometrica*, 21, 503–546.
- Bernoulli, D. (1954). Specimen theoriae novae de mensura sortis [Exposition of a new theory on the measurement of risk]. *Econometrica*, 22, 23–36. (Original work published 1738)
- Blackburn, S. (1998). *Ruling passions*. Oxford, England: Oxford University Press.
- Bolker, E. D. (1967). A simultaneous axiomatization of utility and subjective probability. *Philosophy of Science*, 34, 333–340.
- Broome, J. (1990). Bolker–Jeffrey expected utility theory and axiomatic utilitarianism. *Review of Economic Studies*, 57, 477–502.
- Buchak, L. (2013). *Risk and rationality*. Oxford, England: Oxford University Press.
- Cramer, G. (1975). Letter from Cramer to Nicolas Bernoulli. London, 21 May 1728. In B. L. van der Waerden (Ed.), *Die Werke von Jakob Bernoulli: Vol. 3. Wahrscheinlichkeitsrechnung* [The works of Jakob Bernoulli: Vol. 3. Probability calculus]. Basel, Switzerland: Birkhäuser. (Original work published 1728)
- Davidson, D., McKinsey, J. C. C., & Suppes, P. (1955). Outlines of a formal theory of value, I. *Philosophy of Science*, 22, 140–160.
- Easwaran, K. (2008). Strong and weak expectations. *Mind*, 117, 633–641.
- Egan, A. (2007). Some counterexamples to causal decision theory. *Philosophical Review*, 116, 93–114.
- Gustafsson, J. E. (2013). The irrelevance of the diachronic money-pump argument for acyclicity. *Journal of Philosophy*, 110, 460–464.
- Jeffrey, R. C. (1965). *The logic of decision*. Chicago, IL: University of Chicago Press.
- Jeffrey, R. C. (1983). *The logic of decision* (2nd ed.). Chicago, IL: University of Chicago Press.
- Joyce, J. M. (1999). *The foundations of causal decision theory*. Cambridge, England: Cambridge University Press.
- Joyce, J. M. (2012). Regret and instability in causal decision theory. *Synthese*, 187, 123–145.
- Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decisions under risk. *Econometrica*, 47, 263–291.
- Keynes, J. M. (2000). *A tract on monetary reform*. Amherst, NY: Prometheus Books. (Original work published 1923)

Lauwers, L., & Vallentyne, P. (2016). Decision theory without finite standard expected value. *Economics and Philosophy*, 32(3), 383–407.

McClellenn, E. F. (1990). *Rationality and dynamic choice*. Cambridge, England: Cambridge University Press.

Nover, H., & Hájek, A. (2004). Vexing expectations. *Mind*, 113, 237–249.

Nozick, R. (1969). Newcomb's problem and two principles of choice. In N. Rescher (Ed.), *Essays in honor of Carl G. Hempel: A tribute on the occasion of his sixty-fifth birthday* (pp. 114–146). Dordrecht, Netherlands: Reidel.

Pascal, B. (1967). *Pensées*. New York, NY: Modern Library. (Original work published 1670)

Peterson, M. (2004). From outcomes to acts: A non-standard axiomatization of the expected utility principle. *Journal of Philosophical Logic*, 33, 361–378.

Peterson, M. (2008). *Non-Bayesian decision theory: Beliefs and desires as reasons for action*. Springer.

Peterson, M. (2015). Prospectism and the weak money pump argument. *Theory and Decision*, 78, 451–456.

Peterson, M. (2020). The St. Petersburg paradox. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*. Retrieved from <https://plato.stanford.edu/archives/fall2020/entries/paradox-stpetersburg/>

Peterson, M., & Vallentyne, P. (2018). Self-prediction and self-control. In J. Bermudez (Ed.), *Self-control, decision theory, and rationality* (pp. 48–71). Cambridge, England: Cambridge University Press.

Rabinowicz, W. (1995). To have one's cake and eat it, too: Sequential choices and expected-utility violations. *Journal of Philosophy*, 92, 586–620.

Rabinowicz, W. (2000). Money pump with foresight. In M. J. Almeida (Ed.), *Imperceptible harms and benefits* (pp. 123–154). Dordrecht, Netherlands: Kluwer.

Ramsey, F. P. (1931). Truth and probability (original manuscript, 1926). In R. B. Braithwaite (Ed.), *The foundations of mathematics and other logical essays* (pp. 156–198). London, England: Kegan Paul.

Savage, L. J. (1972). *The foundations of statistics* (2nd ed.). Dover, England: Wiley. (Original work published 1954)

Schick, F. (1986). Dutch bookies and money pumps. *Journal of Philosophy*, 83, 112–119.

von Neumann, J., & Morgenstern, O. (1947). *Theory of games and economic behavior* (2nd ed.). Princeton, NJ: Princeton University Press.

© 2021 The Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

The MIT Press would like to thank the anonymous peer reviewers who provided comments on drafts of this book. The generous work of academic experts is essential for establishing the authority and quality of our publications. We acknowledge with gratitude the contributions of these otherwise uncredited readers.

This book was set in Stone Serif and Stone Sans by Westchester Publishing Services.

Library of Congress Cataloging-in-Publication Data

Names: Knauff, Markus, editor. | Spohn, Wolfgang, editor.

Title: The handbook of rationality / edited by Markus Knauff and Wolfgang Spohn.

Description: Cambridge : The MIT Press, 2021. | Includes bibliographical references and index.

Identifiers: LCCN 2020048455 | ISBN 9780262045070 (hardcover)

Subjects: LCSH: Reasoning (Psychology) | Reason. | Cognitive psychology. | Logic. | Philosophy of mind.

Classification: LCC BF442 .H36 2021 | DDC 153.4/3—dc23

LC record available at <https://lcn.loc.gov/2020048455>