

# Binding and Its Consequences

Christopher J. G. Meacham

May 20, 2009

## Abstract

In “Bayesianism, Infinite Decisions, and Binding”, Arntzenius, Elga and Hawthorne (2004) present cases in which agents who cannot bind themselves are driven by standard decision theory to choose sequences of actions with disastrous consequences. They defend standard decision theory by arguing that if a decision rule leads agents to disaster only when they cannot bind themselves, this should not be taken to be a mark against the decision rule. I show that this claim has surprising implications for a number of other debates in decision theory. I then assess the plausibility of this claim, and suggest that it should be rejected.

## 1 Introduction

In “Bayesianism, Infinite Decisions, and Binding,” Arntzenius, Elga and Hawthorne (2004) examine the significance of *binding*—the ability to irrevocably commit oneself to some future plan of action. They show that in a number of cases, decision theoretic agents who can bind themselves will do much better than decision theoretic agents who cannot. Indeed, in some cases, agents who cannot bind themselves will be driven by decision theory to choose sequences of actions that have disastrous consequences, even when the agents know full well that these choices will lead to disaster, and know ahead of time that these are the choices they’ll make.

One reaction to these cases is to take them at face value, and to conclude that these results are a mark against standard decision theory. If so, this gives us a *prima facie* reason to look more carefully at alternatives which don’t have these consequences, such as the theories of Bratman (1987), Gauthier (1994) and McClennen (1990).

Arntzenius *et al* recommend a different reaction. They suggest that these cases do not give us a reason to be unhappy with standard decision theory. Rather, they argue, the source of these unhappy results is the inability of these agents to bind themselves:

“The lesson is that under certain circumstances, the following ability can be incredibly helpful: ... the ability to irrevocably bind oneself to future courses of action. ... The lack of such ability is not, we say, a deficiency... It’s just that certain situations exploit rational agents who are unable to self-bind.”<sup>1</sup>

We can express this sentiment as follows:

---

<sup>1</sup>Arntzenius, Elga and Hawthorne (2004), p.268–269.

**The Binding Principle:** If a theory of decision making has counterintuitive results that only arise for agents who cannot bind themselves, these results are not a mark against the theory of decision making in question.<sup>2</sup>

If we adopt the Binding Principle, as Arntzenius *et al* suggest, then the cases they examine pose no threat to standard decision theory. Since it is only agents who cannot self-bind who are led to disastrous outcomes, we can attribute these disastrous consequences to their inability to self-bind.

The Binding Principle is interesting for a number of reasons. As we've just seen, it allows standard decision theory to circumvent some otherwise troublesome charges. But it also has consequences for a number of other debates in decision theory. First, it alters the status of the "why ain'cha rich" argument for evidential decision theory. Second, it impacts our assessment of whether decision rules should be self-recommending. Third, it bears on whether decision instability poses a problem for causal decision theory.

Should we adopt the Binding Principle? I'll argue that we should not. I'll suggest that appeals to the Binding Principle mirror earlier appeals to a similar principle regarding mixed acts. And I'll argue that the Binding Principle is problematic for similar reasons.

This paper will proceed as follows. In next section I'll briefly sketch some background. In the third section I'll spell out the implications of adopting the Binding Principle on several debates in decision theory, including the "why ain'cha rich" argument, the question of whether rules should be self-recommending, and decision instability arguments. In the fourth section I'll assess the plausibility of the Binding Principle, and argue that it should be rejected. I conclude in the fifth section by briefly discussing the implications of these verdicts.

## 2 Background: Decision Theory

As I'll understand it, (Bayesian) decision theory can be divided into two parts: a description of the agents to which the theory applies, and a normative claim about how such agents should behave.

The agents to which decision theory applies satisfy the following conditions:

- A1.** The agent's belief state at a time can be represented by a probability function over a space of possibilities. These values, called *credences* or *degrees of belief*, indicate the agent's confidence that the possibility is true, where greater values indicate greater confidence.
- A2.** The agent's evaluative state at a time can be represented by a function which assigns positive real numbers to elements in the space of possibilities. These assignments, called *utilities*, indicate the extent to which the agent values that possibility obtaining, where higher numbers indicate a higher utility.<sup>3,4</sup>

---

<sup>2</sup>This principle is suggested by the discussion in Arntzenius, Elga and Hawthorne (2004), but not explicitly stated. One of the authors has confirmed this understanding of their position in correspondence (though he was clear to note that he was speaking only for himself).

<sup>3</sup>This understanding of utilities sets up decision theory as an account of prudential rationality. Alternatively, one can take these utilities to be whatever is valuable according to some standard, and understand decision theory as an account of instrumental rationality.

<sup>4</sup>If we want to allow for well-defined infinite utilities, then we can use the extended reals to represent utilities instead of the reals.

**A3.** The agent’s potential acts in any decision situation can be represented by a unique set of mutually exclusive propositions  $\{a_1, \dots, a_n\}$  (where  $a_i$  is the proposition that the agent performs the  $i$ th potential act in that situation).

Now consider an agent of this kind who has credences  $cr$  and utilities  $u$ . The *expected utility* ( $EU$ ) for the agent of an act  $a$  is:

$$EU(a) = \sum_{w \in \Omega} cr(w : a) \cdot u_w, \quad (1)$$

where  $\Omega$  is the space of possibilities, and  $cr(w : a)$  is a place holder.<sup>5</sup> By replacing  $cr(w : a)$  with different kinds of functions, (1) yields different kinds of expected utility. If we set  $cr(w : a)$  equal to the agent’s credence in  $w$  conditional on  $a$ , then we get the *evidential expected utility* of the act. If we set  $cr(w : a)$  equal to the agent’s credence that  $w$  would come about were  $a$  to be performed, or something like this, then we get the *causal expected utility* of the act.<sup>6</sup>

The normative part of decision theory claims that agents who satisfy A1-A3 ought to satisfy the following constraint:

**Expected Utility Maximization:** A condition-satisfying agent should only perform a potential act  $a$  if the expected utility of this act is at least as large as the expected utility of any alternatives. I.e., agents should perform acts which maximize expected utility.

By plugging different kinds of expected utility into this constraint, we get different kinds of decision theory. If we plug in evidential expected utility, this constraint yields *evidential decision theory*. If we plug in causal expected utility, we get *causal decision theory*.

Before we proceed, a slight revision to standard decision theory is required. Decision theory is usually understood to prescribe performing one of the acts with the highest expected utility, in the manner just described. But in some cases, there is no such act. For example, suppose an agent will be given  $n$  in dollars, where  $n$  is a natural number of the agent’s choosing. Assume the agent’s utilities are linear in dollars. For any natural number  $n$ , the expected utility of choosing to get  $n$  dollars will increase as  $n$  does. Since there is no largest  $n$ , there is no act with the highest expected utility.

In cases of this kind we can’t expect an agent to choose an act with the highest expected utility, since there is no such act. There are a couple of ways to try to handle this: we might employ satisficing in cases where there’s no highest expected utility act, or we might modify decision theory so that it no longer picks out “best” acts, but instead merely provides a “better than” ordering over them. For the purposes of this paper it will be more convenient to adopt the first approach, so that we can talk about what an agent ought to do, etc., in the usual way.

---

<sup>5</sup>This characterization of (1) assumes a countable number of possibilities, of course. To accommodate uncountably many possibilities, we can extend (1) in the usual way.

<sup>6</sup>I borrow this terminology from Collins (1996). For a discussion of some of the different “causal” ways to understand  $cr(w : a)$ , see Joyce (1999).

## 3 Implications of the Binding Principle

### 3.1 Predictable Disaster

Arntzenius *et al* present a number of cases where the ability to bind oneself seems desirable. Consider the Satan's Apple case:

Satan cuts up an apple into countably infinite pieces in the following way: he cuts the apple in half, and then cuts the remaining half in half, and then cuts the remaining quarter in half, and so on. In a minute he will offer Eve the first piece, 30 seconds later he will offer her the second piece, 15 seconds after that he will offer her the third piece, and so on. So at the end of two minutes he'll have offered her every piece. Now, Eve knows that she'll be expelled from the Garden of Eden—a consequence with -10 units of utility—if she accepts an infinite number of pieces. But Eve likes apples, and the utility she gets from eating a piece is equal to its size ( $1/2$  for eating the first piece,  $1/4$  for eating the second piece, and so on). What should Eve do?

Let's assume that Eve takes her decision about whether to accept each piece to be causally independent of her decisions regarding the other pieces if she's a causal decision theorist, or evidentially independent of her decisions regarding the other pieces if she's an evidential decision theorist.

What does decision theory suggest Eve do with respect to the first piece? Eve believes that taking the first piece won't have any bearing on whether she accepts or declines any of the other pieces. If she's going to accept an infinite number of other pieces, then declining the first piece won't save her from getting expelled from Eden, and accepting the first piece will increase her overall utility by  $1/2$ . If she's only going to accept a finite number of other pieces, then accepting the first piece won't get her expelled from Eden, and eating it will increase her overall utility by  $1/2$ . So at the end of the sequence of offers her overall utility will be greater by  $1/2$  if she accepts the first piece no matter what else she does. So decision theory will tell her to take it.

What about the second piece? The same reasoning applies with respect to the second piece, so decision theory will tell Eve to accept the second piece as well. Likewise, decision theory will tell her to accept every piece she is offered. But if Eve accepts every piece she'll get kicked out of the Garden of Eden, and her overall utility will be 9 units lower than if she had declined every piece.

This looks like a troubling result for standard decision theory. Given standard decision theory, Eve is rationally required to make a series of decisions which she knows will result in disaster: her eviction from the Garden of Eden.

But if we adopt the Binding Principle, as Arntzenius *et al* suggest, then we'll come to a different conclusion. Suppose Eve has the ability to bind herself. When Eve is offered the first piece, she now has an infinite number of actions to choose from: in addition to just accepting or rejecting the first piece, she can opt to bind herself to accept or reject some or all of the other pieces she will be offered in the future. And in this case, decision theory will recommend that she bind herself to accepting some finite number of pieces, and she will end up well-fed and safely ensconced in the Garden.<sup>7</sup>

---

<sup>7</sup>Recall, from section 2, that we're taking decision theory to employ something like satisficing in cases in which there is no best act.

So Eve will only be led to disaster if she lacks the ability to bind herself. Given the Binding Principle, it follows that the disaster that befalls Eve if she lacks the ability to bind herself is not a mark against decision theory. Rather, this case just demonstrates how desirable the ability to bind oneself can be.

### 3.2 “Why Ain’cha Rich?”

In the standard Newcomb’s case, you are presented with the choice of taking the contents of two boxes, or just the contents of the first box. A nearly perfect predictor has attempted to predict your choice. If she thinks you’ll take just the first box, she’ll put a million dollars in it. If she thinks you’ll take both boxes, she’ll leave the first box empty. The second box always contains a thousand dollars.<sup>8</sup>

According to the causal decision theorist, you should always take both boxes. That way, you will be a thousand dollars richer, no matter what the predictor has predicted. According to the evidential decision theorist, you should take only the first box. That’s because the expected monetary reward for choosing one box is higher than that of choosing two boxes.<sup>9</sup>

If we expect the agents who employ one decision making theory to generally be richer than the agents who employ some other decision making theory, this seems to be a *prima facie* reason to favor the first theory over the second. Both causal and evidential decision theorists agree that, in the Newcomb’s case, evidential decision theorists tend to end up wealthier than causal decision theorists. Both expect the evidential decision theorists to get a million dollars when she chooses the first box, and both expect causal decision theorists to get only a thousand when she chooses both boxes. So the Newcomb’s case provides a *prima facie* reason to favor evidential over causal decision theory. As Gibbard and Harper put it, the causal decision theorist faces the question: “if you’re so smart, why ain’t you rich?”<sup>10</sup>

#### 3.2.1 Response 1: Rewarding Irrationality

The standard response to the “why ain’cha rich?” argument is this:<sup>11</sup>

In Newcomb’s case, the predictor will reliably reward “one-boxing”. So those who one-box will reliably end up better off than those who don’t. But this doesn’t show that one-boxing is rational. It merely shows that “if someone is very good at predicting behavior and rewards predicted irrationality richly, then irrationality will be richly rewarded.”<sup>12</sup>

This response shows that the causal decision theorists can provide a consistent explanation for why they don’t take the evidential decision theorist’s wealth to be an indication of rationality. While the evidential decision theorist takes these rewards to be reason to think the pre-rewarded act is rational, the causal decision theorist takes the rewards to be merely a feature of the background situation that is irrelevant to the rationality of the act.

---

<sup>8</sup>See Nozick (1969).

<sup>9</sup>As usual, we’re assuming that the agent’s utilities are linear in dollars.

<sup>10</sup>Gibbard and Harper (1985), p.153.

<sup>11</sup>See Gibbard and Harper (1985), Lewis (1981) and Joyce (1999).

<sup>12</sup>Gibbard and Harper (1985).

But as a response to the “why ain’cha rich?” argument, this isn’t very satisfying. First, this response doesn’t show very much. After all, it’s not surprising that causal decision theory will judge the acts it prescribes to be rational, and those it doesn’t prescribe irrational. This just demonstrates that causal decision theory is consistent. Second, this response won’t cut any ice with the evidential decision theorist, who will maintain that irrational acts *can’t* be predictably pre-rewarded—if the act is predictably pre-rewarded, then it will be the rational act.<sup>13</sup> A more satisfying response to the “why ain’cha rich?” argument would do more than just show that the causal decision theorist’s position is consistent. It would also undermine the evidential decision theorist’s claim that these considerations give us a reason to favor evidential decision theory.

This is where the second line of response to the “why ain’cha rich?” argument comes in. These responses try to do more than just show that causal decision theory is consistent—they also try to undermine the intuition that these considerations provide a *prima facie* reason to favor evidential decision theory. Let’s look at two such “second-line” responses.

### 3.2.2 Response 2: Gibbard and Harper

One response, offered by Gibbard and Harper (1985), attempts to show that the “why ain’cha rich?” arguments can be used against both causal and evidential decision theory in order to support apparently crazy theories of decision making. If so, then the evidential decision theorist should also doubt that we should take “why ain’cha rich?” considerations into account when evaluating theories of decision making. Inessential details aside, the case they consider is this:

As in the standard Newcomb’s case, suppose you must decide between taking one or two boxes, where a predictor has placed a million dollars in the first box iff he predicts you will take one box, and where the second box always has a thousand dollars. In this case, you’ll only be allowed to take the boxes if your decision-making dispositions satisfy certain conditions; but both evidential and causal decision theorists satisfy these conditions.<sup>14</sup> Now suppose that both of the boxes are transparent, so you will see the contents of both boxes before you make your choice. What should you do?

In this case, both causal and evidential decision theory will tell you to take both boxes. And, since the predictor is nearly infallible, both evidential and causal decision theorists will tend to end up with a thousand dollars. By contrast, agents who employ a decision making theory according to which you should only take the first box no matter what you see, will tend to end up with a million dollars. So evidential decision theory seems to be as vulnerable to “why ain’cha rich?” arguments as causal decision theory.

But note that if we adopt the Binding Principle, this response is no longer compelling. An evidential decision theorist who has the ability to self-bind will bind herself to choosing

---

<sup>13</sup>Assuming that this reward outweighs the utility of the alternatives. For more discussion of the evidential decision theorist’s stance on this argument, see Lewis (1981).

<sup>14</sup>You’re not allowed to take boxes if you have a disposition which would make correct prediction impossible. For example, you’re not allowed to take the boxes you choose if your decision making dispositions are: “Take the first box if I see there’s nothing in it, and take both boxes if I see there’s a million in it.” (If your dispositions are such that the predictor can effectively choose which decision you make—you’ll take two boxes if you see nothing in the first box, and just the first box if you see the million—we can assume the predictor is stingy, and won’t put anything in the first box.)

one box before she's shown what's in them. Since the predictor will predict this, she'll put a million into the first box, and the binding evidential decision theorist will end up rich. Since it's only non-binding evidential decision theorists who will end up poor in this case, it follows from the Binding Principle that we shouldn't take this to be a mark against evidential decision theory.

But self-binding causal decision theorists can still end up poor. Consider a version of the Newcomb's case where the predictor makes her prediction before the agent is born. The binding causal decision theorist will be unable to causally influence the prediction, and so she will end up choosing both boxes and getting only a thousand dollars. So even when we restrict our attention to agents who can bind themselves, the "why ain'cha rich" argument against causal decision theory remains.

So if we adopt the Binding Principle, the Gibbard and Harper response is ineffective. The case they discuss is not a problem for the evidential decision theorist, but Newcomb's case is still a problem for the causal decision theorist.

### 3.2.3 Response 3: Arntzenius

Another "second-line" response, offered by Arntzenius (2008), attempts to show that there are cases where we'll expect causal decision theorists to end up richer than evidential decision theorists. If there are such cases, the evidential decision theorist can no longer maintain that evidential decision theorists are generally richer than causal decision theorists, and the "why ain'cha rich" argument for evidential decision theory collapses.

It's not easy to construct such cases. If, for example, we just stipulate that the "evidential decision theory" choice will be punished in some way, then it will no longer be the choice that evidential decision theory recommends (see Lewis (1981)). However, Arntzenius (2008) shows that there are cases in which causal decision theorists will do better than evidential decision theorists. Consider the following case:

Suppose there will be a 10 game series between the Yankees and the Red Sox. You know that the Yankees have a 90% chance of winning any given game. You'll only be allowed to bet on each game if your decision-making dispositions satisfy certain conditions; both evidential and causal decision theorists satisfy these conditions.<sup>15</sup> If you're allowed to bet, then you can bet on either the Yankees (in which case you earn a dollar if the Yankees win, and lose two dollars if the Red Sox win), or the Red Sox (in which case you lose one dollar if the Yankees win, and earn two dollars if the Red Sox win). Finally, before you place each bet, an infallible predictor will tell you whether you'll win or lose the bet. How should you bet?

If no predictor were involved, both causal and evidential decision theory would tell you to bet on the Yankees every time, since the expected utility of betting on the Yankees ( $0.9 \cdot 1 + 0.1 \cdot -2 = 0.7$ ) would be greater than the expected utility of betting on the Red Sox ( $0.9 \cdot -1 + 0.1 \cdot 2 = -0.7$ ). If you bet on the Yankees every time, we'll expect you to win nine times, to lose once, and to be up by \$7 by the end of the series.

---

<sup>15</sup>You're not allowed to bet if you have a disposition which would otherwise make the set-up of the case impossible. For example, you're not allowed to bet if your betting dispositions are: "Bet on the Red Sox if I'm told I'll win my bet, and bet on the Yankees if I'm told I'll lose my bet." Since the predictor can't consistently tell you that you'll win or lose your bet, these dispositions make the set-up of the case impossible.

How should your betting behavior change in light of what the predictor tells you? If you're a causal decision theorist, your behavior won't change at all. The predictor doesn't tell you anything causally relevant to the outcome of the game, so you'll effectively ignore what she says.

If you're an evidential decision theorist, on the other hand, your betting behavior will change. Suppose the predictor tells you that you'll win. Since you get \$2 if you win betting on the Red Sox and only \$1 if you win betting on the Yankees, you'll bet on the Red Sox. If she tells you that you'll lose, you'll also bet on the Red Sox, since you lose \$2 if you lose betting on the Yankees and only \$1 if you lose betting on the Red Sox. So you'll bet on the Red Sox no matter what the predictor tells you. And by doing so, we expect you to win once, lose nine times, and be down by \$7 by the end of the series.

In this case we expect the causal decision theorists to do better than the evidential decision theorists: we'll expect the causal decision theorists to end up \$7 ahead, and we expect the evidential decision theorists to end up down by \$7. So the evidential decision theorist can not maintain that evidential decision theorists are generally better off: they're better off in some situations, but worse off in others. This deflates the "why ain'cha rich?" argument for adopting evidential decision theory instead of causal decision theory.

But note that if we adopt the Binding Principle, matters are different. Suppose the evidential decision theorist has the ability to bind herself. Then in the Red Sox and Yankees case she'll bind herself to betting on the Yankees before the predictor informs her of the outcome of the game. So the binding evidential decision theorist will bet on the Yankees every time, and will end up as rich as the causal decision theorist. Since it's only non-binding evidential decision theorists who will end up poor in the Red Sox and Yankees case, it follows from the Binding Principle that we shouldn't take this to be a mark against evidential decision theory. The argument against the causal decision theorist, however, remains. So if we adopt the Binding Principle, Arntzenius' response to the "why ain'cha rich" argument won't work.

### 3.2.4 Assessing the "Why Ain'cha Rich?" Argument

Without the Binding Principle, the second line of response to the "why ain'cha rich" argument succeeds. The cases that Gibbard and Harper (1985) and Arntzenius (2008) present undermine the evidential decision theorist's claim that "why ain'cha rich" considerations favor evidential decision theory. One might conclude from this that we should ignore "why ain'cha rich" considerations when assessing decision theories. Alternatively, one might conclude that we should be unhappy with both evidential and causal decision theory. But either way, "why ain'cha rich" considerations will fail to support evidential decision theory over causal decision theory.

The final analysis looks different, however, if we adopt the Binding Principle. With the Binding Principle, the second line of response to the "why ain'cha rich" argument fails. While evidential decision theorists who can't bind themselves may end up losing out in the cases Gibbard and Harper (1985) and Arntzenius (2008) present, this isn't a reason to worry about evidential decision theory. Rather, this just demonstrates the unhappy position of agents who are unable to bind themselves. Causal decision theorists can still consistently deny that the evidential decision theorist's position is rational, of course. But they are unable to diffuse the *prima facie* intuition that "why ain'cha rich" considerations are relevant; an intuition which, given the Binding Principle, tells in favor of evidential



decision theory.

### 3.3 Self-Recommendation

Skyrms (1982) raises the question of when decision rules are *self-recommending*:

“The question of what decision method to use for a sequence of decision problems is itself a decision problem. If the rules of rational decision are formulated generally enough, they can be applied to such problems. Let us call a sequence of decision problems a *world*, and the problem of which decision theory to adopt for the individual problems in the sequence, the *world decision problem*. For a given world decision problem, a decision rule might recommend adopting a conflicting rule for dealing with the problems of that world. On the other hand, for certain worlds, certain decision rules will be *self-recommending*.”<sup>16</sup>

More precisely, let a *decision problem* be an ordered triple consisting of the agent’s credences, utilities, and the set of available acts. Let a *comprehensive strategy* be a function from decision problems to acts. Let a *perspective* be an ordered pair consisting of a credence and utility function. Given a decision rule (“*X*-decision theory”) we can work out two things. First, we can work out which comprehensive strategies correspond to the choices an *X*-decision theorist might actually make. Second, we can work out which comprehensive strategies *X*-decision theory takes to be best from a given perspective. So given evidential decision theory, for example, we can work out which comprehensive strategies describe how an evidential decision theorist might act, and we can work out which comprehensive strategies evidential decision theory takes to be the best from a given perspective—the comprehensive strategies which get assigned the highest evidential expected utility. A decision rule is *self-recommending* from a perspective when the comprehensive strategies that describe the behavior of a rule-following agent are among the strategies that the rule considers best. So a decision rule is self-recommending from a perspective when it takes itself to be a good decision rule to adopt.

Being self-recommending is a nice feature for a decision rule to have. A decision making rule that is self-recommending is robustly confident about the acts it prescribes.

That said, we shouldn’t expect a decision rule to be self-recommending from every perspective. For instance, we shouldn’t expect a decision rule to be self-recommending from the perspective of an agent who believes her future credences will fail to cohere with her current ones in certain ways. Consider an agent who knows that a given coin toss landed heads, and knows that she will forget this information by the time she’s in a position to bet on it. Given this, evidential decision theory may assign a higher evidential expected utility to comprehensive strategies that describe an agent who always bets on heads than to strategies that describe how an evidential decision theorist would act. But this doesn’t indicate that evidential decision theory lacks confidence in its prescriptions. Rather, evidential decision theory doesn’t endorse the comprehensive strategies that describe evidential decision theorists because after a certain point the evidential decision theorists will be making choices with faulty credences.

Likewise, we shouldn’t expect a decision rule to be self-recommending from the perspective of an agent who believes her future utilities will differ from her current ones. Consider, for example, an agent who initially only values the happiness of sentient beings,

---

<sup>16</sup>Skyrms (1982), p.707.

but believes she will come to value only money at some point in the future. Given this, evidential decision theory will assign a higher evidential expected utility to comprehensive strategies which describe an agent who spends her life promoting happiness than to comprehensive strategies that describe the evidential decision theorist, who will soon turn to collecting money at the expense of others. But again, this doesn't indicate that evidential decision theory fails to be confident in its own prescriptions. Rather, evidential decision theory doesn't endorse the comprehensive strategies that describe the evidential decision theorist because after a certain point the evidential decision theorist will be making choices with different utilities.

Finally, we shouldn't expect a decision rule to be self-recommending from the perspective of an agent who thinks she'll come to believe (rightly or wrongly) that she may deviate from the prescriptions of the rule. Consider an alcoholic who enjoys the atmosphere at the local bar, but who believes that she will give into temptation and start drinking if she goes there. According to evidential decision theory, the alcoholic should not go to the bar. But evidential decision theory assigns a higher evidential expected utility to comprehensive strategies which prescribe going to the bar and not drinking than to strategies which prescribe staying home. Again, this doesn't indicate that evidential decision theory lacks confidence in its prescriptions. Rather, these two ways of picking out comprehensive strategies diverge because the first way factors in the possibility of failing to adhere to a strategy, while the second way—assessing how good adhering to a strategy would be—does not.

To avoid these kinds of cases, let's restrict our attention to the perspectives of agents who believe that they will (a) update by conditionalization, (b) have static utilities, and (c) have a negligible credence that they'll deviate from the decision rule. Given these restrictions, are evidential and causal decision theory self-recommending?

No. To see that evidential decision theory can fail to be self-recommending, recall the Gibbard and Harper (1985) variant of Newcomb's case described in the previous section, where you see the contents of both boxes before you make your choice. In this case evidential decision theory will recommend that you take both boxes, regardless of what you see. And since the predictor will predict this, evidential decision theorists will tend to only get a thousand dollars. Now consider the decision rule *X*-decision theory, which prescribes taking the first box no matter what you see. Since the predictor will predict this, *X*-decision theorists will tend to get a million dollars. So the evidential expected utility of acting in accordance with *X*-decision theory will be higher than the evidential expected utility of acting in accordance with evidential decision theory.

To see that causal decision theory can fail to be self-recommending, consider the Satan's Apple case. The causal expected utility of acting like a causal decision theorist will be low—she'll be expelled from Eden. On the other hand, the causal expected utility of acting in accordance with a decision rule that results in her taking only the first 100 pieces will be much higher—she'll remain in Eden and get most of the apple.

So neither evidential nor causal decision theory are self-recommending in all of the cases we'd like. This is a *prima facie* reason to look for some other decision rule.

But if we adopt the Binding Principle, we'll come to a different assessment. If we further restrict our attention to agents who can bind themselves, evidential and causal decision theory will both be self-recommending. An evidential decision theorist who can self-bind in the Gibbard and Harper case will bind herself to choose only the first box before she sees their contents. And the evidential expected utility of these comprehensive strategies

will be as high as the evidential expected utility of the comprehensive strategies prescribed by any other decision rule. Likewise, a causal decision theorist who can self-bind in the Satan's Apple case will bind herself to accepting only a finite number of pieces, and the causal expected utility of this kind of comprehensive strategy will be as high as the causal expected utility of the comprehensive strategies prescribed by any other decision rule. So given the Binding Principle, both evidential and causal decision theory are appropriately self-recommending.

### 3.4 Decision Instability

One of the worries that has been raised for causal decision theory is that it leads to cases of *decision instability*. We have a case of decision instability if, for every available act  $A$ , the expected utility of  $A$  conditional on  $A$  is lower than the expected utility of some other act  $A'$  conditional on  $A$ . In such cases there's a sense in which you'll be displeased with your choice no matter what, since as soon as you choose an act you'll believe some other act is better.<sup>17</sup>

A classic example of decision instability is the Death in Damascus case presented by Gibbard and Harper (1985):

“Suppose the man knows the following. Death works from an appointment book which states time and place; a person dies if, and only if, the book correctly states in what city he will be at the stated time. The book is made up weeks in advance on the basis of a highly reliable prediction. An appointment on the next day has been inscribed for him. Suppose, on this basis, the man would take his being in Damascus the next day as strong evidence that his appointment with Death is in Damascus, and would take his being in Aleppo the next day as strong evidence that his appointment is in Aleppo.”<sup>18</sup>

If the man's only choices are to go to Aleppo or to go to Damascus, what should he do?

According to causal decision theory, the man's decision is unstable. Conditional on the assumption that he'll go to Aleppo, the causal expected utility of going to Damascus will be higher, since he'll expect Death to be at Aleppo. Conditional on the assumption that he'll go to Damascus, the causal expected utility of going to Aleppo will be higher, since he'll expect Death to be at Damascus. So no matter what the man decides to do, as soon as he makes his choice he'll come to believe that the other act is better.

Given evidential decision theory, on the other hand, decision instability cannot arise. Since you calculate the evidential expected utility of an act by assessing the relevant probabilities conditional on the act, the evidential expected utility of going to Aleppo and the evidential expected utility of going to Aleppo conditional on going to Aleppo will be the same.<sup>19</sup> (Indeed, the evidential expected utility of going to Damascus conditional on going to Aleppo is not well defined, so the means of comparing the two acts that decision

---

<sup>17</sup>Assuming that you update by conditionalization, and that the evidence you get from performing an act is just that you've performed the act.

<sup>18</sup>Gibbard and Harper (1985), p.154–155.

<sup>19</sup>Though there are variants of evidential decision theory which are also subject to decision instability. I ignore these variants, since I am only concerned with the canonical version of evidential decision theory here.

instability requires is unavailable.)<sup>20</sup>

Is this kind of case a problem for causal decision theory? A number of people have thought so, and have proposed revisions of causal decision theory in order to avoid this kind of instability.<sup>21</sup> That said, it's not clear that the decision instability that appears in cases like Death in Damascus is sufficient to justify these revisions. First, as Gibbard and Harper note, it's not clear that this result is counterintuitive. Arguably, decision instability is exactly what we should expect in this kind of case:

“Any reason the doomed man has for thinking he will go to Aleppo is a reason for thinking he would live longer if he stayed in Damascus, and any reason he has for thinking he will stay in Damascus is a reason for thinking he would live longer if he stayed in Aleppo. Thinking he will do one is a reason for doing the other. That there can be cases of unstable [causal expected utility]-maximization seems strange, but the strangeness lies in the cases, not in [causal expected utility]-maximization: instability of rational decision seems to be a genuine feature of such cases.”<sup>22</sup>

Second, the decision instability in this case doesn't result in any strange behavior: the man will still go to either Aleppo or Damascus. So the fact that his decision is unstable doesn't seem particularly troubling.<sup>23</sup>

But more troubling worries for causal decision theory arise when we consider multiple-act versions of the cases where decision instability arises. For example, consider the following variant of the Death in Damascus case, inspired by Richter (1986):

Add to the Death and Damascus case the following details. The man is halfway between Damascus and Aleppo. He must decide whether to take a step toward Damascus, or a step toward Aleppo. After he's taken his first step, he has to decide whether to take his second step toward Damascus or Aleppo, and so on. It will take him an hour to reach either city, and there are 5 hours left before nightfall. He knows that if he is not in a city by nightfall, Death's jackal companions will come and eat him, and this will be a much more unpleasant way to die than meeting Death in the city. Finally, he's thirsty, tired and hungry, and he would like to get to a city as soon as possible. What should the man do?

---

<sup>20</sup>We can adopt primitive conditional probabilities to get around this obstacle, of course. But decision instability will still not arise for evidential decision theory for the reason just given.

<sup>21</sup>For example, see Sobel (1983), Eells (1985), and Weirich (1985).

<sup>22</sup>Gibbard and Harper (1985), p.156.

<sup>23</sup>Egan (2007) presents some other cases in which decision instability arises, and argues that in these cases causal decision theory delivers the wrong verdicts. Arntzenius (2008) suggests that the proponent of causal decision theory can reasonably deny that the verdicts in question are counterintuitive. But in any case, Egan's concern is not with decision instability *per se*, but with the fact that he thinks the causal decision theory verdicts are counterintuitive.

Weirich (1985) and Richter (1986) raise a different kind of worry: they argue that in the kinds of cases in which decision instability arises for causal decision theory, it should generally be rationally permissible to know what you're going to do before you do it. But one cannot both satisfy causal expected utility-maximization in cases of decision instability and know what you're going to do before you do it. So if we grant that knowing what you're going to do before you do it is rationally permissible in cases like the Death in Damascus case, causal decision theory is in trouble.

If the man acts in accordance with causal decision theory, he will act as follows. He will start going toward one city until his credence that Death will be there starts to increase. Then he'll start taking steps back toward the other city until his credence shifts again. And so on.<sup>24</sup> He'll continue to dither about unhappily in the desert until he has no more time left to spare if he's to make it to one of the cities before nightfall. At that point it will become a choice between going to the closer city or be eaten by jackals, and so he'll head toward the nearest city.

This seems like an unhappy result for causal decision theory. If he adheres to causal decision theory, the man can know ahead of time that he'll spend 4 extra hours dithering about unhappily in the desert instead of relaxing and sipping margaritas at a local tavern, and yet he'll do it anyway.

But if we adopt the Binding Principle, we should view these cases in a different light. If the man has the option of binding himself, then he won't dither about in the desert—he'll choose one of the two cities, and bind himself to taking steps in that direction the rest of the way. So agents who can bind themselves won't suffer from the kind of problem posed in the Dithering in Damascus case. Thus, given the Binding Principle, we shouldn't take these cases to indicate problems with causal decision theory. Rather, we should take these cases to be further demonstrations of the fact that agents who are unable to bind themselves can end up in unhappy situations.

## 4 Assessing The Binding Principle

Adopting the Binding Principle is an appealing way of maintaining the *status quo* with respect to standard decision theory. But there are reasons to doubt it. To make these doubts salient, let's first look at a similar issue that arises in discussions about decision instability.

### 4.1 Decision Instability and Mixed Acts

Causal decision theory has been criticized for giving rise to decision instability. One response to this worry has been to note that the instability tends to go away if we allow for *mixed acts*.<sup>25</sup> A mixed act is typically thought of as a decision to base one's act on the result of some chance event. For example, one might flip a fair coin in order to decide between performing act *A* or act *B*. More precisely, call the acts typically considered in decision theory—pull a lever, accept a bet, etc.—*pure acts*. We can characterize a mixed

---

<sup>24</sup>One might worry about whether his credence *should* increase. After all, there's no reason to think he can't predict how he's likely to behave ahead of time. And why should his credence that he will end up at Aleppo increase as he takes more steps toward Aleppo if he knows he's going to change his mind and turn around?

Of course, similar reasoning can be applied if we claim that his credence that he'll end up at Aleppo should remain the same. If his credence in ending up at Aleppo will remain the same as he takes steps toward Aleppo, then we'd expect him to end up at Aleppo, in which case it seems his credence that he'd end up at Aleppo should have been increasing after all.

In any case, we can side-step the issue by stipulating that, after every step, the man has a chance of involuntarily taking another step in a random direction (and after that random step, a chance of taking yet another step in a random direction, and so on). With this addition, the man's credence that he'll end up at Aleppo *should* increase after he takes a step toward Aleppo, based solely on these chances.

<sup>25</sup>Harper (1986).

act as a probability distribution over these pure acts, where these probabilities are independent of the possible outcomes of the pure acts in question. In the example just given, the mixed act assigns a probability of 50% to *A* and a probability of 50% to *B*.

Responses to the decision instability that employ mixed acts have generally been met with skepticism.<sup>26</sup> One such response is to claim that agents like us do have access to mixed acts, and thus won't face cases of decision instability. But this claim is implausible: we don't generally have access to the probability-generating devices required, and it's hard to see how we could perform mixed acts without them. Another response is to restrict the application of decision theory to agents who have access to mixed acts, and *a fortiori* restrict decision theory to cases where decision instability won't arise. But this robs standard decision theory of much of its interest. We want to know how agents like ourselves should act, and if standard decision theory doesn't provide this guidance, we'll have to look elsewhere.

More generally, it seems like shifting to cases in which the agents have access to mixed acts is changing the topic. After all, the mixed acts version of a decision problem is a different decision problem than the original. In the original Death in Damascus case, for example, the agent has only two choices: go to Damascus or go to Aleppo. But if we take the agent to also have mixed acts available then the agent has an infinite number of choices: the two pure acts of the original case, and the mixed acts consisting of every probability distribution over those pure acts.

Moreover, the addition of these choices substantially alters the nature of the case. Not only are we assuming that the agent has an adjustable chancy device and the willpower to commit himself to certain acts given certain outcomes, we're also assuming that the outcomes of the chance device are independent of the outcomes of the pure acts in question. This is a substantial assumption. In order to maintain this independence in the Death and Damascus case, we have to assume that Death, astounding predictor that he is, isn't able to predict the outcomes of the agent's chance device. This makes it clear that the agent in the mixed acts version of the case is in a very different situation from the agent in the original case. In the original case, the agent has virtually no chance of staying alive, since Death is a fantastically good predictor. In the mixed acts variant the agent is in a much better position: since Death can't predict the outcome of her chance device, if she uses it to randomly determine which city she'll go to, her chance of surviving will be 50%.

A different way of using mixed acts to ward off decision instability problems is this:

**The Mixed Acts Principle:** If a theory of decision making has counterintuitive results that only arise for agents without access to mixed acts, these results are not a reason to reject the theory of decision making in question.

This response grants that agents like us might face cases of decision instability on causal decision theory. But since it's only agents who don't have access to mixed acts who will face such cases, they should not be taken to be marks against causal decision theory. Rather, the moral is that having access to mixed acts is desirable. Agents without such access might run into uncomfortable situations, like cases of decision instability. But this doesn't indicate a flaw in the decision rule in question; rather, it's just a demonstration of why it's nice to have access to mixed acts.

It's true that having access to mixed acts is desirable. In the Death in Damascus case, for example, the agent who has access to mixed acts has a much better chance of staying

---

<sup>26</sup>See Arntzenius (2008) for criticisms of this kind.

alive. But this isn't enough to justify the Mixed Acts Principle. Being omniscient is desirable too, but that doesn't mean we can ignore any counterintuitive consequences a decision rule has for agents who are not.

I take the Mixed Acts Principle to be implausible. The mixed acts version of a case is simply a different case from the original. If a decision rule delivers a counterintuitive result in the original case, then that seems to be a mark against the theory. And whether this consequence also appears when we consider the mixed acts version of the case is irrelevant to our evaluation of the original case.

The Mixed Acts Principle may derive some apparent plausibility from a related claim:

**Ought entails Can (Mixed Acts):** If no decision rule can avoid a given counterintuitive results without employing mixed acts, then it's not a demerit of a particular decision rule that it can't avoid these results without employing mixed acts.

This claim is plausible. But unlike the Mixed Acts Principle, it does not support the mixed acts response to the decision instability worries since there are decision rules which won't lead to decision instability, even when mixed acts are not available. As we've seen, evidential decision theory is such a rule.<sup>27</sup> So this claim doesn't allow the causal decision theorist to ignore cases of decision instability that arise for agents who don't have access to mixed acts.

Causal decision theory might still be the most plausible decision rule, of course. One's assessment of a theory depends on lots of factors, of which this is only one. But none of that changes the fact that decision instability seems to be a demerit of the theory. And when we assess the pros and cons of the theory, it should be taken as such.

## 4.2 The Binding Principle

We've seen some reasons to be skeptical of the Mixed Acts Principle. These same considerations apply to the Binding Principle.

The binding version of a decision problem, like the mixed acts version, is a different decision problem than the original: the agent is presented with a number of additional acts to choose from. And like the Mixed Acts Principle, the Binding Principle is *prima facie* implausible. If a decision rule delivers a counterintuitive result in a given case, then that seems to be a mark against the theory. And whether this counterintuitive result also appears when we consider some other case—a case where the agent has binding acts available, say—is irrelevant to our evaluation of this case.

To put it another way, if we accuse standard decision theory of allowing agents who can't bind themselves to choose acts which lead to disaster, the proponent of the Binding Principle will respond: "Well, if they had access to self-binding they'd be alright." But this is just to say: "Well, if they were in some other situation instead of this one then they'd be alright." This is true, of course, but how does it make the fact that agents are led to disaster in this situation any less problematic?

There is a claim that resembles the Binding Principle that may lend it plausibility:

**Ought entails Can (Binding):** If no decision rule can avoid a given counterintuitive result without employing binding, then it's not a demerit of a particular decision rule that it can't avoid these results without employing binding.

---

<sup>27</sup>There are, of course, many others. For example, any decision rule in which one's credences don't play a substantive role.

This claim is plausible. But unlike the Binding Principle, it does not lend itself to the defense of standard decision theory, since there *are* decision rules which won't lead agents who can't bind themselves to disaster.

For instance, we can construct a rule which effectively advises an agent to choose the acts that she would have wanted to bind herself to. Let  $iu(\cdot)$  represent the agents initial utilities, let  $ic(\cdot)$  represent the agents initial credences, and let  $CS$  be the proposition that the agent will act in accordance with a particular comprehensive strategy. Let the “committal expected utility” of a comprehensive strategy be:

$$CoEU(CS) = \sum_i ic(w_i : CS) \cdot iu(w_i). \quad (2)$$

As before, we can get causal and evidential versions of this theory by substituting the appropriate replacement in for  $ic(w_i : CS)$ . Now consider a decision rule which requires agents to satisfy the following constraint:

**Committal Expected Utility Maximization:** A condition-satisfying agent should perform the act picked out by a comprehensive strategy which maximizes committal expected utility.

This “committal decision theory” prescribes the acts that the agent would have bound herself to had she the ability to do so.<sup>28</sup> And non-binding agents who adopt committal decision theory will not be led to disaster in the kinds of cases Arntzenius *et al* describe.

It's worth noting that committal decision theory doesn't require agents to have willpower, plans or foresight. Nor does it require agents to have actually committed themselves to future courses of action at some earlier time.<sup>29</sup> Like standard decision theory, committal decision theory is simply a rule which prescribes acts to agents in decision problems. And it requires no more of agents than standard decision theory does.

Likewise, it's worth noting that committal decision theory *does* take what you've learned into account when prescribing acts. Committal decision theory prescribes the acts selected by the comprehensive strategies which maximize committal expected utility. And comprehensive strategies are functions from decision problems—ordered triples consisting of the agent's credences, utilities and the set of available acts—to a subset of the available acts. So even though committal decision theory is insensitive to your current credences when evaluating comprehensive strategies, the comprehensive strategies themselves are sensitive to your current credences when recommending acts. As a result, what you've learned does end up getting taken into account by committal decision theory, once we get down to the level of which acts you ought to perform.<sup>30</sup>

---

<sup>28</sup>The exact content of this theory depends on how we choose to understand the  $ic$  function (2) employs. One possibility is to take it to take it at face value, as the agent's first credence function. Another possibility is to take it to be something like the subject's “ur-priors”—the credences the agent ought to have if she had no evidence whatsoever. (This ur-prior would be the same for all agents for objective Bayesians, and different for different agents for subjective Bayesians.) A related possibility is to take  $ic$  to be something like the initial credences of an ideal subject in the agent's situation. This would allow us to think of committal decision theory as a kind of “ideal observer theory” of prudential rationality. (Thanks to Dennis Whitcomb for suggesting this third possibility.)

<sup>29</sup>In this respect, committal decision theory differs from the proposals offered by Bratman (1987), Gauthier (1994) and McClennen (1990).

<sup>30</sup>A number of people have expressed something like the following worry: “Why should the agent choose acts that seem reasonable according to her initial credences? Shouldn't she choose acts that seem reasonable according to her current credences, not her initial ones?” This is a natural worry. But it is difficult to spell it out



Committal decision theory is one rule which avoids the counterintuitive consequences that Arntzenius *et al* describe, but there are many others. Similar rules have been proposed by Bratman (1987), Gauthier (1994) and McClennen (1990). And other rules with similar features are not hard to find.

Indeed, we can even find rules which, like standard decision theory, are functions of the agent’s *current* credences and utilities. For example, starting with (2), we can replace the agent’s initial utilities with her current ones. And we can replace the agent’s initial credences with the expectation of what the agent currently believes her initial credences might have been. Call the resulting expression the *committal expected utility\** of a comprehensive strategy:

$$CoEU^*(CS) = \sum_i cr(ic_i) \sum_j ic_i(w_j : CS) \cdot u(w_j). \quad (3)$$

We can plug committal expected utility\* into the constraint described above to get *committal decision theory\**. And, as before, non-binding agents who adopt committal decision theory\* will escape disaster in the kinds of cases Arntzenius *et al* describe.

## 5 Conclusion

The Binding Principle that Arntzenius *et al* employ has a number of interesting consequences. It bolsters evidential decision theory by undercutting the second line of response to the “why ain’cha rich” argument. It bolsters causal decision theory by undercutting the decision instability arguments. It allows standard decision theory to escape the blame for the disastrous outcomes it leads agents to in the Satan’s Apple case. And it provides standard decision theory with an excuse for failing to be self-recommending in certain ways. All told, the Binding Principle offers an appealing way of maintaining the front-runner status of standard decision theory.

Unfortunately, the Binding Principle is implausible. The Binding Principle is analogous to the Mixed Acts Principle, and is problematic for the same reasons. As a result, the various marks against standard decision theory described above still apply.

This makes rules like (2) and (3) more attractive by comparison. These rules escape “why ain’cha rich” and decision instability arguments, avoid disasters in scenarios like the Satan’s Apple case, and are appropriately self-recommending. For similar reasons, the proposals of Bratman (1987), Gauthier (1994) and McClennen (1990) are more attractive in this light.

---

in a compelling way.

One might be asking why an agent should do what’s reasonable according to her initial credences instead of what’s reasonable according to her current credences. But if we’re assuming that what is reasonable is what standard decision theory prescribes—expected utility maximization—then this is question begging. And if we’re assuming that what is reasonable is what committal decision theory prescribes, then an agent who satisfies committal decision theory *is* doing what’s reasonable according to her current credences.

Alternatively, one might be asking why one ought to believe that committal decision theory is the right decision rule. But the answer to this question is straightforward: we’re justified in thinking that committal decision theory is the right rule to the extent to which it provides the intuitively correct prescriptions. And, as we’ve seen, there are several cases in which committal decision theory arguably provides more intuitive prescriptions than standard decision theory.

Of course, some form of standard decision theory may still be the most plausible decision rule. One's assessment of a rule depends on lots of factors besides those considered here. And rules like (2) and (3) have counterintuitive consequences of their own.<sup>31</sup> But none of this changes the fact that the problems described above are demerits of standard decision theory. And when we assess the pros and cons of standard decision theory, they should be taken as such.<sup>32</sup>

## References

- Arntzenius, Frank. 2008. "No Regrets, or: Edith Piaf Revamps Decision Theory." *Erkenntnis* 68:277–297.
- Arntzenius, Frank, Adam Elga and John Hawthorne. 2004. "Bayesianism, Infinite Decisions, and Binding." *Mind* 113:251–283.
- Bratman, Michael. 1987. *Intentions, Plans and Practical Reason*. Harvard University Press.
- Collins, John. 1996. "Supposition and Choice: Why 'Causal Decision Theory' is a Misnomer." Presented at the CUNY Graduate Center Philosophy Colloquium.
- Eells, Ellery. 1985. "Weirich on Decision Instability." *Australasian Journal of Philosophy* 63:473–478.
- Egan, Andy. 2007. "Some Counterexamples to Causal Decision Theory." *Philosophical Review* 116:93–114.
- Gauthier, David. 1994. "Assure and Threaten." *Ethics* 104:690–716.
- Gibbard, Allan and William Harper. 1985. Counterfactuals and Two Kinds of Expected Utility. In *Paradoxes of Rationality and Cooperation: Prisoner's Dilemma and Newcomb's Problem*, ed. R Campbell and L Sowden. University of British Columbia Press.
- Harper, William. 1986. "Mixed Strategies and Ratifiability in Causal Decision Theory." *Erkenntnis* 24:25–36.
- Joyce, James. 1999. *The Foundations of Causal Decision Theory*. Cambridge University Press.
- Lewis, David. 1981. "Why ain' cha rich?" *Nous* 15:377–380.
- McClennen, Edward. 1990. *Rationality and Dynamic Choice*. Cambridge University Press.

---

<sup>31</sup>For example, the evidential versions of these rules will prescribe the one-boxing response to the Gibbard and Harper case discussed in section 3.2.2. And these rules will recommend that the alcoholic discussed in section 3.3 go to the bar, even though she believes she will probably start drinking if she does.

<sup>32</sup>I would like thank Frank Arntzenius, Philip Bricker, Maya Eddon and Dennis Whitcomb for helpful comments and discussion.

- Nozick, Robert. 1969. Newcomb's Problem and Two Principles of Choice. In *Essays in Honor of Carl G. Hempel*, ed. Nicholas Rescher. Reidel: Dordrecht.
- Richter, Reed. 1986. "Further Comments on Decision Instability." *Australasian Journal of Philosophy* 64:345–349.
- Skyrms, Brian. 1982. "Causal Decision Theory." *The Journal of Philosophy* 79:695–711.
- Sobel, Howard. 1983. "Expected Utilities, and Rational Actions and Choices." *Theoria* 49:159–183.
- Weirich, Paul. 1985. "Decision Instability." *Australasian Journal of Philosophy* 63:465–472.