

Desire, Expectation and Invariance

March 15, 2015

Abstract

The Desire-as-Belief thesis (DAB) states that a rational person desires a proposition exactly to the degree that she believes or expects the proposition to be good. Many people take David Lewis to have shown the thesis to be inconsistent with Bayesian decision theory. However, as we show, Lewis' argument was based on an Invariance condition that itself is inconsistent with the (standard formulation of the) version of Bayesian decision theory that he assumed in his arguments against DAB. The aim of this paper is to explore what impact the rejection of Invariance has on the DAB thesis. Without assuming Invariance, we first refute all versions of DAB that entail that there are only two levels of goodness. We next consider two theses according to which rational desires are intimately connected to expectations of (multi-levelled) goodness, and show that these are consistent with Bayesian decision theory as long as we assume that the contents of 'value propositions' are not fixed. We explain why this conclusion is independently plausible, and show how to construct such propositions.

1 Introduction

The Desire-as-Belief thesis (DAB) holds that a rational person desires a proposition exactly to the degree that she believes the proposition to be desirable or good. David Lewis, the originator of the thesis, considered it to be a version of anti-Humeanism about motivation, since, if the thesis is true, then forming beliefs about what is good would suffice to produce the requisite desires and hence to motivate action. Lewis and others have also claimed that the thesis is important for both cognitivism and objectivism in

ethics, respectively the views that moral judgements are beliefs and that these beliefs can be true.

How exactly Desire-as-Belief relates to these meta-ethical views is, however, in need of clarification. As we point out in the first section below, Humeans and anti-Humeans, and objectivists and subjectivists, might all want some version of the Desire-as-Belief thesis to be true. So what is at stake, in the debate over the thesis, is not these meta-ethical views. Rather, the issue is whether rationality requires a particular relationship between an agent's desire for a proposition *A* and her beliefs about propositions that express the value of *A*, where this value can be interpreted in line with all of the above mentioned meta-ethical views.

Lewis (1988) famously argued that the Desire-as-Belief thesis must be false, since it conflicts with Bayesian decision theory.¹ In section 3 we reproduce Lewis' argument and show that it was based on a false assumption, namely Invariance, according to which the desirability of a proposition is independent of its probability and truth-value. As we will explain, this assumption is not only intuitively implausible but inconsistent with the (standard formulation of the) version of Bayesian decision theory that Lewis assumed in his arguments against DAB.

The main aim of this paper is to explore what consequences the rejection of Invariance has on the debate over the DAB thesis. In section 4 we provide a counterexample to (both causal and evidential versions of) the most simple formulations of DAB, according to which a rational agent desires a proposition to the extent that she *believes* the proposition to be good or desirable. According to somewhat more sophisticated versions of DAB, rational agents' desires are intimately connected to their *expectation* of goodness. In section 5 we consider one such thesis, called *Desire-as-Expectation*, which both Humeans and anti-Humeans about motivation might want to accept, and which is not undermined by our counterexample. We show that despite appearances, the thesis is perfectly consistent with Bayesian decision theory. We conclude (in section 6) with a comparison between Desire-as-Expectation and a generalised version of DAB that allows for multiple levels of goodness, and show that our argument for Desire-as-Expectation also serves

¹John Collins (1988) similarly showed that a qualitative version of DAB is inconsistent with plausible constraints on qualitative belief revision. We will be concerned only with quantitative versions of DAB in this paper.

to vindicate the generalised Desire-as-Belief thesis.

2 Why be interested in the Desire-as-Belief thesis?

In this section we explain why the Desire-as-Belief thesis is important. Lewis himself took the thesis to be a version of anti-Humeanism about motivation. So although he recognised that there might be other ways of formulating the anti-Humean view—and hence that a refutation of DAB was not sufficient to show anti-Humeanism to be false—he did think that the truth of DAB would entail the falsity of the Humean view on motivation. Moreover, he thought that DAB entailed ethical objectivism, i.e., the view that there are objective truths about ethical reality. We believe that Lewis was mistaken on both accounts. As we show below, Humeans can be committed to the Desire-as-Belief thesis (section 2.1) and so might ethical subjectivists be (section 2.2). So the debate around the DAB is not one about the truth of these meta-ethical views. Instead, the question is what rationality can consistently require of agents capable of forming attitudes to propositions about value.

2.1 DAB is not anti-Humean

The Humean view on motivation that Lewis endorsed states that “we are moved entirely by desire: we are disposed to do what will serve our desires according to our beliefs.” (1988: 323) His aim in refuting the Desire-as-Belief thesis was, as previously mentioned, to refute one version of the opposing anti-Humean view. The “Anti-Humean’s main thesis”, Lewis says, is that there are necessary connections between people’s desires and their beliefs about what is good: “It is just impossible to have a belief about what would be good and lack the corresponding desire.” (ibid: 324). Hence, if someone truly believes that it would be good to help an old lady cross the street, then he would necessarily desire to do so, and would thus (Lewis thinks the anti-Humean would say) be motivated to help the lady cross the street.

Having explained how he understood the distinction between the Humean and anti-Humean view about motivation, Lewis formulated the Desire-as-Belief thesis—which,

recall, states that a rational agent desires a proposition to the extent that she believes the proposition to be good—as one version of anti-Humeanism. However, it should be evident that DAB is not at all anti-Humean. The thesis does not say that people in general necessarily desire what they believe to be good, nor that these are one and the same mental attitude. What it does say, is that *rational* agents maintain a particular relationship between their desires and their beliefs about the good. This is at least how we will be interpreting the thesis, and, indeed, how the thesis must be understood for it to be at all interesting (as Ruth Weintraub (2007) points out). For if DAB were a psychological claim about ordinary people, then we wouldn't need a philosophical or decision-theoretic argument to examine its plausibility: citing ordinary psychological experience (with all its akrasia and confused desires) would then suffice to refute the thesis.

Given that DAB is a claim about rationality, many Humeans about motivation will happily accept it. In a well-known defence of the Humean theory, Michael Smith (1994) says: “Humeans ... need not deny the contingent coexistence of beliefs and desires ... nor ... that the contingent coexistence of certain beliefs and desires is rationally required” (1994: 119; see also Smith (1987)). But that is exactly what the DAB thesis (as we understand it) states: it does not claim that any beliefs are *necessarily* accompanied by some particular desires, nor that some beliefs are identical to some desires. So it is not a claim about the constitution of desires and beliefs, nor a thesis about what motivates ordinary people. Instead, it is a thesis about *rationality*.

The above shows that even if DAB is true, Humeanism about motivation need not be false. What if DAB generally fails to hold: does that mean anti-Humeanism must be false as well? No: there might still be *some* fixed connection between peoples' desires and their beliefs about the good; just not the precise relationship postulated by DAB. (As we explain in section 5, this relationship might be captured by John Broome's thesis of *Desire-as-Expectation*, even if the DAB thesis is false.)

2.2 DAB is not an objectivist thesis

In his second paper on the Desire-as-Belief thesis, Lewis claimed that the truth of DAB would vindicate objectivism in ethics. If there are value propositions, belief in which are connected to desires in the way postulated by DAB, then “some of them presumably would be true”, Lewis says, and then “we surely would want to say that the true ones were objective truths about ethical reality.” (1994: 307)

Contrary to Lewis’ claim, it does not follow, from the existence of value propositions, that some of them are true. They might, for instance, all be indeterminate, as Graham Oddie (2001) points out. Or they might all be false. For it is possible that all propositions of the form *A is good* are false, and so are all propositions of the form *A is not good*. Value propositions would then share a peculiar feature with sentences such as ‘the present king of France is bald’ (assuming that Bertrand Russell (1905) was right in his interpretation of such sentences). That is, propositions of the form *it is not the case that A is good* are true, and so are propositions of the form *it is not the case that A is not good*; but all propositions of the form *A is good* and *A is not good* are false. It would, admittedly, have rather strange implications if the DAB thesis were true while all value propositions are false: it would entail that a rational agent who knows the truth about value should not desire anything. But the point is that the existence of value propositions does not, by itself, entail that some of them are true.

More importantly, there are variants of the Desire-as-Belief thesis that will seem plausible to many ethical subjectivists. The literature on the DAB thesis has mostly interpreted the value propositions that figure in the formal statement of the thesis as expressing the claim that some proposition is objectively good. However, as will become apparent in next section when we give a precise statement of DAB, these value propositions can just as well be interpreted as expressing that some proposition is desirable to some agent, or simply that it satisfies her desires. Now consider the requirement that a rational agent desires a proposition *A* to the extent that she believes *A* will satisfy her considered desires. Many subjectivists would be happy to accept this version of DAB as a requirement of rationality. Most, if not all of us often violate this requirement since we often desire things while knowing that they won’t satisfy us. But subjectivists might

want to say that this is never true of an ideally *rational* person. As we shall see, Lewis' argument, if successful, also refutes this subjectivist version of DAB.

An argument by Oddie (1994) could nevertheless be seen as showing that the falsity of DAB would be especially unwelcome news for cognitivists and other objectivists in ethics. Suppose evaluative judgements are beliefs, as cognitivists claim. Then an ideal agent would, Oddie suggests, maintain *harmony* between her desires and her ethical beliefs, and thus desire whatever she believes to be morally good or right. In other words, the ideal agent would satisfy DAB. If in addition ethical objectivism is true, then the moral beliefs of an ideal agent would all be true, Oddie suggests, and she would, moreover, only desire that which is in fact good. So again, an ideal agent would satisfy DAB. However, if Lewis' arguments against DAB are successful, then an ideal agent cannot be both harmonious and fully rational. So although Lewis' argument doesn't show that either objectivism or cognitivism is false, it does create a bit of a dilemma for cognitivists and objectivists, since they must accept that ideal agents are either irrational or inharmonious.

The above dilemma also arises for many types of subjectivists however. For instance, we can, as previously mentioned, interpret the value propositions that figure in the formal statement of DAB as expressing that some proposition satisfies the desires of some agent. Therefore, if Lewis' argument against DAB is successful, then subjectivists must accept that even the most rational agents cannot maintain harmony between their desires and their beliefs about what satisfies these desires. Hence, the falsity of DAB would be no less awkward for subjectivists than it is for objectivists.

2.3 DAB is a thesis about rationality

What is at stake in the debate over Desire-as-Belief is not the status of different theories of value and human motivation. Rather, the issue is whether it is plausible that *rationality* requires that there be a fixed quantitative relationship between an agent's desire that *A* and her beliefs concerning propositions about the value of *A*, where this value can be understood in a way that is acceptable to the proponents of each of the views on meta-ethics and human motivation that we have discussed above.

On the face of it, it does seem that there should be some such relationship. Given that we can understand ‘desirable’ in a way that best suits our theory of morality and human nature, many subjectivists and objectivists alike, and also both Humeans and anti-Humeans about motivation, would, we think, want to say that a rational agent desires a proposition to the degree that she believes it to be desirable. But if Lewis’ argument against DAB is correct, then that cannot be true, since it follows from it that a rational agent cannot satisfy DAB. If, instead, we want to say that an *ideal* agent desires a proposition to the degree that she believes it to be desirable, then we are forced to conclude, from Lewis’ arguments, that the ideal agent is not rational. So, to sum up, a complete rejection of DAB-like theses would have implications for a much broader class of theories of value and human motivation than Lewis seems to have realised.

3 Lewis’ argument against DAB is unsound

In this section we examine Lewis’ argument against the Desire-as-Belief thesis and show that it relies on an assumption called *Invariance* (section 3.1), before explaining why we think this assumption is false and hence Lewis’ argument unsound (section 3.2).

3.1 Lewis’ argument requires Invariance

To state Lewis’ argument against DAB, let A, B, \dots be propositional variables, understood as sets of possible worlds, and \mathcal{W} be the set of all worlds. Ω is the set of all propositions—i.e., the set of all subsets of \mathcal{W} — P a subjective probability (or credence) function from Ω into the interval $[0,1]$ and V a desirability function. P thus measures the degrees of belief of a rational agent and V the strength of her desires.² Let P_A be an agent’s revised credence function after learning that A , which, if the agent revises her beliefs in accordance with Bayesian conditioning, is equal to $P(\cdot | A)$. Likewise let V_A be the agent’s revised desirability function after she has learned that A . Finally, from any proposition A , we construct the *halo-proposition* \mathring{A} , interpreted as a proposition about the

²[X] has pointed out that although functions like V are standardly called ‘desirability’ functions—following Jeffrey’s (1983)—it might be more appropriate to call them ‘desiredness’ functions, since they need not represent ‘actual’ (or objective) desirability. Nevertheless, we will stick with the convention, and call V a desirability function.

value of A ; for instance, that A is desirable or that A is good.

Lewis made a number of arguments against DAB (Lewis (1988), Lewis (1996)). We will state the simplest of these, found in section 4 of his second paper on ‘Desire as Belief’ and for which it must be assumed that the desirability measure V is bounded by 0 and 1. First, a formal statement of the version of DAB he considers:

Thesis 1 (Desire-as-Belief (DAB)). *For any A and according to any rational agent:*

$$V(A) = P(\mathring{A}) \quad (1)$$

In what follows we will, from time to time, refer to this version of DAB as the *simple* DAB thesis, to distinguish it from the more complicated version that we discuss in sections 5 and 6.

Lewis also assumed *Invariance*, according to which the desirability of a proposition is invariant under changes in its probability, and hence, in particular, unaffected by whether the proposition is true or false:^{3,4}

Assumption (Invariance (INV)). *For any A and according to any rational agent:*

$$V_A(A) = V(A) \quad (2)$$

Together DAB and Invariance imply that $P(\mathring{A}) = V(A) = V_A(A) = P_A(\mathring{A})$. (The last equality holds since DAB is assumed to continue to hold after a rational agent learns that A .) In other words, A and \mathring{A} are probabilistically independent:

Implication (Independence (IND)). *For any A and according to any rational agent:*

$$P(\mathring{A}) = P_A(\mathring{A}) = P(\mathring{A} | A) \quad (3)$$

Why is IND problematic? It is not hard to show that even if we start with a probability function for which such independence holds, it is not guaranteed that it will continue

³Lewis’ first argument against DAB contained an Invariance assumption that was limited to maximally specific propositions (see e.g. Lewis (1988): 327). As should be apparent in section 3.2, Jeffrey’s decision theory, within which Lewis’ discussion of DAB takes place, is also inconsistent with that version of the Invariance assumption (given how Jeffrey understands the desirability of the tautology).

⁴Costa, Collins and Levi’s (1995) argument against DAB also relies on Invariance.

to hold after the agent in question revises her beliefs in accordance with Bayesian conditionalisation (an example is given in next paragraph). That is, suppose that an agent's revised partial beliefs after she has learned A , represented by the probability function P_A , is related to her partial beliefs *before* learning A , represented by P , by the following condition: for any proposition B , $P_A(B) = P(B | A) = P(A \wedge B)/P(A)$. Then we cannot be sure that this agent will satisfy IND both before and after such revision. Hence, given Invariance, a person cannot satisfy DAB unless she fails to update via Bayesian conditionalisation (i.e., unless she violates what we'll call BAYES). So if, as Lewis assumed, Invariance is true and BAYES is a requirement of rationality, then DAB cannot be rationally required (assuming that rationality does not make inconsistent demands).

Here is an example where IND and BAYES cannot both be satisfied. Assume that there is some proposition A such that $0 < P(A), P(\bar{A}) < 1$. (If we cannot make these assumptions, without undermining DAB, the thesis only holds in quite trivial cases, as Lewis points out.) This of course implies that $0 < P(A \vee \bar{A})$. Hence, it should be possible for an agent to learn that $A \vee \bar{A}$ and we should have no problems with conditionalising on this proposition, using Bayesian conditioning. Moreover, suppose IND holds before the agent in question learns $A \vee \bar{A}$; i.e., $0 < P(\bar{A} | A) = P(\bar{A}) < 1$. Now the problem is that given these assumptions, when P is updated by $A \vee \bar{A}$ using Bayesian conditionalisation, IND no longer holds: This update leaves the conditional probability of \bar{A} given A unchanged, but increases the probability of A (since $P(A) < 1$). Hence, when $P(\bar{A}) = P(\bar{A} | A)$ and $0 < P(A), P(\bar{A}) < 1$, $P_{A \vee \bar{A}}(\bar{A}) \neq P_{A \vee \bar{A}}(\bar{A} | A)$. Less formally, if IND holds before an agent has learned $A \vee \bar{A}$, then she cannot still satisfy IND after having learned this unless she violates BAYES.

3.2 Why Invariance is false

In making his argument against the Desire-as-Belief thesis, David Lewis drew on the version of decision theory developed by Richard Jeffrey (1983). According to Jeffrey's theory, the desirability of any proposition, A , is a weighted average of the different mutually exclusive and jointly exhaustive possibilities compatible with A , where the

weight on each of these possibilities, S_i , is given by $P(S_i | A)$. More formally:

Desirability (Jeffrey's formula) For any $A \in \Omega$ such that $P(A) > 0$ and any partition $\{S_i\}$ of \mathcal{W} :

$$V(A) = \sum_{S_i \in \mathcal{W}} P(S_i | A) \cdot V(A \wedge S_i) \quad (4)$$

The standard interpretation of this measure, suggested to Jeffrey by Leonard Savage (Jeffrey (1983): 82), is that it measures the *news value* of a proposition. In other words, the V -value of a proposition represents how much the agent in question would welcome the news of its truth. But it can be interpreted more generally as measuring the value of the difference that the proposition's truth makes, relative to the agent's expectations.

Since propositions are taken to be sets of possible worlds, there is just one tautological proposition, denoted by T , which is the set of all possible worlds. It follows from Jeffrey's formula that:

$$V(T) = V(B \vee \neg B) = V(B) \cdot P(B) + V(\neg B) \cdot P(\neg B) \quad (5)$$

Jeffrey assumed that the desirability of the tautology was a constant; conventionally zero, which denotes the 'neutral' point in the desirability scale. This is entirely natural on the news-value interpretation of desirability as the news that the tautology is true is really no news at all. Indeed, on any interpretation of desirability in terms of the value of the difference that the truth of a proposition makes, the tautology will be neither desirable nor undesirable (hence, 'neutral') since its being true makes no difference at all (given that its truth is always already given).⁵ Lewis' assumption that the desirability function is bounded by zero and one however rules out this conventional choice of zero for the tautology.⁶ But the zero-normalisation is not essential here. What is important is that because the tautology is an 'empty' proposition whose truth is always certain, its desirability does not depend on the truth or falsity of any other proposition. Hence its desirability does not change as a result of learning the truth of any contingent proposition.

⁵For example, a similar argument can be made in terms of the willingness-to-give-up interpretation of desirability: If you are certain, both before and after learning B , that T is true already, then you should be willing to give up the same—i.e., nothing of any value—to make T true before and after learning B .

⁶Unless all contingent propositions are equally desirable.

To see why this implies the falsity of Invariance, consider what happens as the agent's probability for some proposition B rises. As $P(B)$ approaches 1, $P(\neg B)$ approaches 0, and therefore $V(\neg B).P(\neg B)$ approaches 0. Hence, since $V(B).P(B) + V(\neg B).P(\neg B) = V(T)$, $V(B)$ must approach $V(T)$. In other words, as B becomes more probable its desirability approaches the desirability of the tautology; so "one who believes that a proposition is true cannot desire that it be true" (Jeffrey (1983): 63). But we cannot, of course, assume that all propositions are always considered equally desirable as the tautology. Hence, the desirability of a proposition, according to an agent, is generally not independent of her credence in the proposition (contra Invariance).

As Jeffrey notes, the idea that a desire for A is 'neutralised' when one comes to believe that A is true, had been defended well before he wrote *The Logic of Decision*. In Plato's (360 BC) *Symposium*, for instance, Socrates claims that it is constitutive of desire that one cannot desire that which one already has:

[E]very one who desires, desires that which he has not already, and which is future and not present, and which he has not, and is not, and of which he is in want ...

If propositions are the objects of desires, as both Lewis and Jeffrey assumed, then Socrates can be read as claiming that he who desires, desires that which he does not already believe to be true. Socrates is of course not claiming that a person should desire *not* having that which he already has. Instead, the idea seems to be that by acquiring something the desire for that thing becomes neutral, just as Jeffrey's theory entails.

These implications of Jeffrey's framework accord well with intuition about news-value. If you are almost certain that you will survive the day, then you won't be very excited to learn that you will indeed survive the day; that won't be much news to you. However, you would presumably be devastated to learn that it is not true. This is in agreement with what Jeffrey's measure entails: As a desirable proposition becomes more and more probable, the desirability of its truth approaches the neutral point in the desirability measure 'from above', while the desirability of its negation generally moves further and further down the negative part of the desirability scale.

It also accords well with other ‘difference-making’ theories of value. In economics, for instance, it is commonly held that one can measure the extent to which a person desires a proposition (or good) by what the person would be willing to give up to make that proposition come true (or to acquire that good). On this understanding of how to measure the strength of a person’s desires, it is clear that degrees of desire do not satisfy Invariance. How much a person is willing to give up in order to make a proposition A come true, is certainly not independent of whether she takes A to be true already. Suppose the person considers A , whose truth she is uncertain about, to be desirable. In other words, she would be willing to give up at least something of value in order to make A true; let’s call this something G . What about when she considers what she would be willing to give up to make A true, after having learned that A is true already? Surely, whatever she would be willing to give up to make A true after having learned that A is true already (if anything at all), should be less valuable to her than G .

3.3 Lewis’ Argument for Invariance

Invariance, we have argued, is incompatible with standard interpretations of Jeffrey’s concept of desirability. This is not enough to dismiss Lewis’ argument against DAB, however, because Lewis evidently took the anti-Humean to be committed to Invariance. Indeed, his argument for Invariance suggests that he thought that the very formulation of the DAB thesis entailed a commitment to it. If this were true, our argument against Invariance would best be interpreted as a diagnosis of *why* the DAB thesis conflicts with decision theory rather than a refutation of Lewis’ argument.

But why did Lewis think that an anti-Humean was committed to Invariance? In arguing for this he noted that Invariance holds for any proposition if it holds for all the maximally specific subcases making it up. But if a subcase,

were maximally specific merely in all ‘factual’ aspects, . . . , then it would be no surprise if a change in belief changed our minds about how good it would be [that the subcase were true] . . . But the subcase was supposed to be maximally specific in *all* relevant aspects . . . The subcase has a maximally specific hypothesis about what would be good built right into it. So in

assigning it a value, we do not need to consult our opinions about what is good. We just follow the built-in hypothesis.

(Example. How good would it be if, first, pain were the sole good, and second, we were all about to be in excruciating and everlasting pain?—I have to say that this would be good, and so I value the case highly. My opinion that in fact pain is no good does not affect my valuing of the hypothetical case in which, *ex hypothesi*, pain is good. My opinion does cause me to give the case negligible credence, of course, but that is different from affecting the value.) (Lewis (1988): 332)

Lewis' argument can be broken into two distinct claims. First, that a maximally specific subcase entails the truth or falsity of all propositions about goodness. And second, that the desirability that a rational agent may assign to such a subcase is uniquely determined by the content of such propositions. In the simple case under consideration, for instance, the only relevant goodness propositions are the halo propositions and goodness comes in only two degrees; zero or one. So Lewis' two claims amount in this case to the assertion that there exist maximally specific propositions that entail, for any proposition A , either that \hat{A} or that $\neg\hat{A}$ and that for any such proposition A , $V(A \mid \hat{A}) = 1$.

Lewis meant his claims to extend to the more general case in which there are multiple levels of goodness. We discuss these more fully in later sections, but for the moment let us simply accept the existence of propositions that express the fact that the truth of A is good to degree i and denote them by \hat{A}_i . Then what Lewis requires more generally is the truth of what might be called the Principal Moral Principle; the claim that the desirability that A is true, conditional on A 's truth being good to degree i , is just i .⁷ More formally:

Principal Moral Principle: For any A and according to any rational agent:

$$V(A \mid \hat{A}_i) = i$$

Lewis' argument is seductive, but misleading. Let us grant that the anti-Humean should

⁷See Nissan-Rozen (2015) for a discussion of this principle.

accept both the existence of goodness propositions and the truth of the Principal Moral Principle. It does not follow, however, that they must accept Invariance. For Lewis' argument requires not just that there are goodness propositions but that, for *every* proposition, there exists a corresponding halo proposition, including those propositions that are maximally specific both with respect to the non-evaluative facts and the goodness facts. But then it follows that there must exist what one might call *self-evaluative* propositions: propositions that make claims about their own goodness. In particular, there must exist a maximally specific proposition, MAX, which is of the form 'A is the case, B is not the case, . . . , and it would be good if MAX were the case'.

It is the self-evaluative maximally specific propositions that impose the dubious Invariance condition on Jeffrey's framework. Not by any means the first instance in philosophy of self-referential propositions causing trouble! Lewis seems to think that the anti-Humean must accept the existence of such self-evaluative propositions. But this is quite implausible. Acceptance of the existence of goodness propositions such as 'it would be good if it were to rain tomorrow' does not force acceptance of second-order goodness propositions such as 'it would be good if it were good that it would rain tomorrow', let alone propositions that are maximally specific with respect to higher-order goodness claims.

Nor does formulation of the DAB thesis require the acceptance of higher order goodness propositions. For instance, Bradley and List (2009) exhibit a framework in which the DAB thesis can be satisfied but in which self-evaluative propositions do not exist. The key idea is to distinguish a set \mathcal{W} of factual worlds, subsets of which are the purely factual propositions, from a set \mathcal{V} of evaluative worlds, subsets of which are purely evaluative propositions, such as the halo-propositions or goodness-level propositions. Intuitively, we can think of factual worlds as capturing all the physical facts and evaluative worlds as capturing all the goodness facts about these physical facts. A basic (or maximally specific) possibility in this framework is just a pair (w, v) such that $w \in \mathcal{W}$ and $v \in \mathcal{V}$, and an extended proposition is any set of such world-pairs.

Now notice that on this construction there exists, for every factual proposition A , a corresponding halo proposition \mathring{A} , but no haloed halo propositions. Nonetheless, the

framework is rich enough to state a version of the simple DAB thesis adequate for both those Humeans and those anti-Humeans that want to endorse it. In particular, since the extended propositions form a Boolean algebra we can define a Jeffrey desirability function V and a probability function P on the set of them and then meaningfully require that for any *factual* proposition A , $V(A) = P(A)$. Invariance, on the other hand, is simply not required by the framework. There is thus a sense in which Lewis' argument for Invariance rests on a careless formal construction, one which allows for indiscriminate 'haloing' of propositions!

3.4 Concluding Remarks on Invariance

We have shown both that Invariance is at odds with standard interpretations of desirability and that the DAB thesis does not require it. Our view is that it should therefore be dispensed with. But some of those interested in the DAB thesis might nonetheless be uncomfortable with abandoning Invariance. Ethical objectivists in particular might consider it unacceptable for rational belief about the goodness of a proposition to depend on how probable it is that the proposition is true. Hence, as both Weintraub (2007) and Daskal (2010) indirectly point out, there might seem to be a need to consider other possible responses to the tension between Invariance and standard interpretations of (Jeffrey) desirability. But this perception is based on a misunderstanding. Invariance fails in Jeffrey's framework because of the way desirability is cardinalised and in particular because it is normalised with respect to the tautology. But this does not imply that the goodness *ordering* of worlds varies with changes in belief. On the contrary, if we take the ordering of worlds to be objective, then we can construe desirability as a normalised measure of goodness, that coheres with the betterness ranking in the sense that for all worlds w and w' with non-zero probability, $V(w) \geq V(w')$ just in case w ranks at least as high as w' in the objective betterness ranking of worlds. So an ethical objectivist can subscribe to a DAB thesis formulated within Jeffrey's framework, despite the fact that it implies that rational belief about the goodness of some proposition A varies with changes in the probability of A , because this variation is a feature of the quantitative representation of desirability, rather than a reflection of any changes in the underlying

betterness ordering.

This argument will not move the kind of objectivist who takes numerical desirabilities to be primitive, rather than representations of an underlying betterness ordering.⁸ But the idea that desirability numbers are primitive seems very implausible to us, and is certainly something that most decision theorists would reject. And objectivists should not accept this idea either, at least if they take desirability (or value) to be similar to the quantities, such as length or temperature, that we find in the natural sciences.⁹ For it is a widely held contention in the theory of measurement that such quantities are not primitive, but simply representations of comparative relations such as 'longer than' or 'warmer than'. In any case, it is incumbent on the objectivist who does take numerical desirabilities to be primitive to explain where they come from and how they are measured.

For those not persuaded by these arguments, there are two other possible responses to the inconsistency between Invariance and standard interpretations of Jeffrey's framework. First, one might try to detach the DAB thesis from Jeffrey's decision theory. And second, one could retain Jeffrey's framework but give it a non-standard interpretation that makes it reasonable for the desirability of the tautology to vary. In the rest of this paper, however, we take Invariance to be false and explore the implications for the debate over the DAB thesis.

4 A new counterexample to Desire-as-Belief

Lewis' argument against the DAB thesis has been shown to fail because it is based on an unsound premise. This does not of course mean that the thesis is true. Indeed, in this section, we present an example that refutes all versions of what we above called the *simple* Desire-as-Belief thesis, which is the version of the thesis that has received almost all the attention in the literature on DAB, and according to which an agent desires a proposition to the extent that she *believes* the proposition to be good. The example highlights the implausibility of assuming that goodness comes only in two degrees, as

⁸Thanks to [Y] for bringing this view to our attention.

⁹We thank [Z] for suggesting this analogy.

the simple DAB thesis entails. We first explain the example and show how it undermines a simple *evidential* DAB thesis (section 4.1). A number of authors have proposed a *causal* decision-theoretic version of this simple DAB thesis (see, for instance, Byrne and Hájek (1997) and Williams (2010)). But as we show in section 4.2, our counterexample also undermines a simple causal DAB thesis. The conclusion of this section is that a minimal requirement on DAB theses is that they allow for multiple degrees of goodness. In sections 5 and 6 we consider such versions of the thesis.

4.1 Counterexample to simple evidential DAB

Suppose we are sailing with our two good friends, Ann and Bob, when suddenly both of them fall overboard and find themselves in an equally difficult situation and threatened with drowning. In that situation we *fully* believe the proposition that it would be good that Ann is saved. (Call this proposition \mathring{A} .) Or if we can only fully believe a tautology, then we at least believe \mathring{A} as (or almost as) strongly as we believe any contingent proposition. Nothing in our example hangs on treating Ann and Bob equally, but to simplify the discussion, let us suppose that our feelings for the two are identical in all relevant respects. Thus we also believe the proposition that it would be good to save Bob (call this proposition \mathring{B}) as (or almost as) strongly as we believe any contingent proposition. So for us, in that situation, $P(\mathring{A}) = P(\mathring{B})$ is close to 1. To make the discussion that follows more precise, let us assume that $P(\mathring{A}) = P(\mathring{B}) = 1 - \gamma$.

In the situation we are imagining, we would also *desire* very strongly that Ann is saved (proposition A), and would desire equally strongly that Bob is saved (proposition B).¹⁰ But we find it much more desirable—in fact about twice as desirable—that *both* Ann and Bob are saved than that only one of them is. So assuming that the status quo (i.e., what happens without an intervention) is that both of them drown, and, moreover, that the probability that one of them is saved is independent of the probability that the other is, we find that $V(A \wedge B)$ is roughly twice $V(A)$. But then the Desire-as-Belief thesis dictates that the probability that it is good that both Ann and Bob are saved, $P(A \wedge B)$, should be close to twice $P(\mathring{A})$. But since $P(\mathring{A})$ is close to 1, $P(A \wedge B)$ can never be close

¹⁰Assuming that the probability of them being saved is roughly equal. The significance of this assumption will become clear at the end.

to twice $P(\bar{A})$). Thus assuming that the requirements of rationality never ban what is rationally permissible, and if we take the attitudes towards Ann and Bob expressed in the above example to be rationally permissible, it seems that DAB cannot be a requirement of rationality.

To save DAB a proponent of it could argue that, contrary to appearances, the attitudes towards Ann and Bob assumed in the counterexample are in fact irrational. There are three ways she could do this. First, she can deny that $P(\bar{A})$ is rationally permitted to be close to 1. Second, she can argue that $V(A) = V(B)$ should be no greater than $(1 - \gamma)/2$. Third, she can argue that $V(A \wedge B)$ should not be much greater than $V(A) = V(B)$. Alternatively, she could argue in a quite different vein, that our assumptions do not have a clear meaning.

Let us take each response in turn. Since we are assuming that saving both Ann and Bob is close to twice as desirable as saving one of them, the first response only works if we require that $P(\bar{A})$ is less than 0.5. So for this response to work, we must be less certain in the proposition that it is good to save Ann (or Bob) than the proposition that a fair coin lands heads up when tossed. It is highly implausible that this is a rationality requirement on beliefs.

In fact, we can make things much worse. Suppose now that we are sailing with not just two but a number of our dear friends when suddenly all of them fall overboard. For the first response to the above counterexample to work, the credence we assign the proposition that it is good to save any one of our friends must get smaller and smaller as we increase the number of people that we imagine to have fallen overboard. But no matter how many friends we have, and how many of them we take out sailing, we would always be almost certain that it would be good to save each of them after having fallen overboard. To take an example, suppose we are sailing with six friends who all fall overboard. Then we cannot be more certain in the proposition that it would be good to save any particular friend than in the proposition that a dice shows side six when rolled! A conception of rationality that requires this seems very implausible.

The second route to saving the DAB thesis involves requiring that $V(A) \leq (1 - \gamma)/2$. But however we interpret desirability, it is hard to believe that rationality requires that

saving Ann (or Bob) be confined to the bottom half of the desirability scale. In any case, this requirement coupled with the Desire-as-Belief thesis implies that $P(\mathcal{A}) \leq (1 - \gamma)/2$. In other words, this response requires us to be less certain in the proposition that it would be good to save Ann than in the proposition that a fair coin comes heads up if tossed. So this second response to our counterexample in the end comes down to the same as the first response and is no less implausible. And again, we can make this response even less plausible by increasing the number of people we are imagining to be in the water.

The third response consists in denying that it is permissible to judge it much more desirable to save both Ann and Bob than just one of them. Ordinary intuition (and many welfarist theories) suggests that saving both Ann and Bob would be roughly twice as desirable as saving one of them. Saving both might be *more* than twice as desirable as saving just one of them; for instance if we feel guilt for choosing to save one of them over the other, or if choosing to save one over the other creates some sort of injustice or unfairness. Or it might be slightly *less*, for instance if Bob and Ann hate each other and would be happier if the other were dead. But if we set these complementarities aside then we are left with the core judgement upon which the example is based: that the desirability of saving Ann (or Bob) is independent of whether the other is saved or not. But if this is so then it would seem to follow immediately from the assumption that saving Bob is equally desirable as saving Ann, that saving both is twice as desirable as saving one.

Could there be complementarities that we are rationally required to give weight to and which make the assumed judgement irrational? It is hard to imagine what they could be. But even if there are such complementarities they are unlikely to make enough of a difference. The problem is that to rescue DAB from our counterexample, it would need to be demonstrated that the difference between the desirability of saving Ann and of saving both must of rational necessity be very small. Suppose, for instance, that we are 90% sure that it would be good to save Ann and 95% sure that it would be good to save both Ann and Bob, so that by DAB, $V(A) = 0.9$ and $V(A \wedge B) = 0.95$. Then it is just slightly above 5% more desirable to save both Ann and Bob than to save one of them.

That is implausible. Having saved one of our friends, we would still make a great deal of effort and be willing to risk or pay quite a lot to save the other. And it is hard to see why that would be irrational. We could of course argue about the plausibility of the exact numbers, but so long as the difference in the probability of \emptyset and $A \wedge B$ is not great, the difference in desirability between saving both friends and just one of them must be *very* small for this last response to work; much smaller than what most people would intuitively accept.

So let us consider the final possible response to the counterexample, which works by questioning the meaningfulness of the assumption that the desirability of saving Ann and Bob is twice that of saving Ann. In the decision-theoretic framework in which DAB is stated, desirability functions are just numerical representations of preferences and only those properties of desirabilities that are analogues of properties of preferences should be considered meaningful. But the notion of ‘twice as desirable as’ fails this test, as is evidenced by the fact that a linear transformation of a desirability function (and in particular one based on different choice of zero point) will yield another desirability function that serves equally well to represent the underlying preference relation, but does not preserve properties such as one prospect being twice as desirable as another. So the counterexample trades on an unsustainable interpretation of desirabilities.

This objection is half-correct. It is true that a linear transformation of a desirability function does not preserve the property that we are interested in. But such transformations are ruled out by the simply version of the DAB thesis, which itself forces a particular choice of the zero and unit scaling points on the desirability function (namely, the certainly bad and the certainly good propositions). One may well object that such a choice of scale is arbitrary, but this would be a reason to object to (this version of) DAB directly. Here we assume for the purposes of the argument that the scaling of desirabilities enforced by the thesis is acceptable and then show that it leads to unacceptable conclusions.

What then does ‘twice as desirable as’ mean within the scope of desirabilities as regulated by the DAB thesis? Roughly this: An agent who regards prospect X as twice as desirable as prospect Y is one who is indifferent between Y being true for certain and

a lottery which makes X true with probability one-half and the certainly bad prospect true otherwise. Suppose for instance that it is certainly bad that both Ann and Bob are not saved. Then it is twice as desirable that both Ann and Bob are saved as that Ann is saved, just in case the prospect of Ann being saved is just as desirable as the prospect of either both being saved or neither, with an equal probability of each.¹¹

To sum up: The attitudes towards our friends Ann and Bob expressed in the above example are rationally permissible, and any attempt to save the simple Desire-as-Belief thesis in light of this counterexample forces us to have attitudes that seem counterintuitive and are certainly not rationally required. Hence, the example shows that this version of DAB must be false.

4.2 Counterexample to simple causal DAB

A number of authors have argued that a version of the Desire-as-Belief thesis that is formulated in terms of causal decision theory rather than Jeffrey's evidential decision theory can withstand Lewis' criticism (see, for instance, Byrne and Hájek (1997), Oddie (2001) and Williams (2010)). As we now show, however, our counterexample undermines all simple causal versions of DAB.

The main difference between causal and evidential decision theory is that the former weights consequences by probability under subjunctive supposition, or the probability of a counterfactual, where the latter weights them by conditional probability. More precisely, let $P_A^\square(w_i)$ measure the probability that world w_i *would be* the case if A *were* true. Then causal decision theory prescribes maximisation of the *causal efficacy value*, U , of an 'action proposition', which is given by:¹²

¹¹It might be objected that this definition assumes that the agent is risk neutral. But the objection is misplaced. Lewis formulated DAB within the decision theory developed by Richard Jeffrey (1983). And in Jeffrey's framework—and, indeed, in all standard decision theories—risk attitudes are built into the desirabilities of propositions in the sense that the method for constructing a cardinal measure of desirability assumes risk neutrality with respect to desirability. Agents are not, however, assumed to be risk neutral with respect to specific goods, and agents who are, say, risk averse with respect to a particular good are modelled with a desirability function that is concave over that good.

¹²Some causal decision theorists (for instance Lewis (1981)) are happy to use Jeffrey's formula for desirability, but suggest we use this causal-efficacy formula for choice-worthiness. Others disagree and argue that desirability should match choice-worthiness (see e.g. Byrne and Hájek (1997)). This disagreement is irrelevant to the present discussion, since for causal decision theory to save DAB, it has to be the case that whatever we call the type of value that figures in the DAB thesis, it is formalised by equation 6.

$$U(A) = \sum P_A^\square(w_i).V(w_i) \quad (6)$$

Correspondingly the version of DAB that the aforementioned authors propose states that:

Thesis 2 (Simple Causal DAB). *For any A and according to any rational agent:*

$$U(A) = P(\mathcal{A}) \quad (7)$$

Although these authors have not explicitly taken issue with Invariance, it is worth noting that this assumption is clearly not valid for causal efficacy value. This is perhaps best illustrated by the *Newcomb* decision problem (Nozick (1969)) that historically provided the main motivation for causal decision theory. Recall that in this decision problem, taking both boxes is evidence for, but does not cause, the emptiness of the ‘black’ box that could contain the larger amount of money. Now let A be the proposition that the agent takes both boxes. On the assumption that A , the black box is (almost certainly) empty. Hence, the utility of A given A is (very close to) the utility of receiving only what is in the ‘opaque’ box that can only contain the smaller amount of money. However, if we make the standard assumptions that causal decision theorists make when suggesting two-boxing, then the (unconditional) utility of A is far greater than the utility of receiving only what is in the opaque box. So Invariance fails for the utility measure that figures in causal decision theory, and Lewis’ argument against the DAB thesis does not work against a causal version of the thesis.

But let us now see how the Simple Causal DAB (SCDAB) fares in light of the counterexample we discussed above. Any causal decision theorist would, we contend, say that the causal efficacy of $A \wedge B$ is roughly twice that of A only. To put it in the terminology that a causal decision theorist is most likely to relate to: The consequence of the act of successfully saving both Ann and Bob is roughly twice as valuable as the consequence of the act of successfully saving only Ann. However, a causal decision theorist will, just like anyone else, presumably be almost certain that it would be good to save Ann. But for the reasons discussed above, the above two judgements cannot

both be true, if SCDAB is correct: If $U(A \wedge B)$ is roughly twice $U(A)$, then by SCDAB, $P(A \wedge B)$ is close to twice $P(A)$, which is inconsistent with the judgement that $P(A)$ is close to one. So our counterexample undermines a simple causal version of DAB.

5 Desire-as-Expectation

In reaction to David Lewis' criticism of DAB, John Broome (1991) proposed the *Desire-as-Expectation* thesis (DAE), which avoids the criticisms that we (and Lewis) have directed against DAB. To state his thesis formally, let $\{G_i\}$ be a partition of the set of possible worlds according to how good they are, such that, for instance, *goodness-level proposition* G_j expresses the fact that the world is good to degree j . Then Broome's thesis says:

Thesis 3 (Desire-as-Expectation). *For any A and according to any rational agent:*

$$V(A) = \sum_i i.P(G_i \mid A) \quad (8)$$

Broome claims that DAE is more plausible as an anti-Humean view than the one Lewis formulated. For there is no reason anti-Humeans should take there to be an equality (or identity) between desires and beliefs; instead, they should simply say that certain desires *result* from beliefs alone. And if we assume (as anti-Humeans should) that the G_i partition is determined by the beliefs of the agent we are modelling, then according to DAE, $V(A)$ is determined by the evaluative beliefs of that agent.

Moreover, Broome thinks the DAE thesis should be no less acceptable to Humeans than anti-Humeans:

Both groups can agree that one should desire something to a degree equal to the expectation of good from it. Where they differ is over what ultimately determines the goodness of a world. A Humean thinks goodness must ultimately be determined by people's desires; an Anti-Humean thinks this is not so. (Broome (1991): 265)

In other words, while anti-Humeans think the G_i -partition is determined by the agent's beliefs, Humeans take the partition to be determined by the agent's desires.

The fact that DAE allows for multiple degrees of goodness makes it deal well with the counterexample we raised in last section to the simple DAB thesis. Without having to specify the candidate levels of goodness, it seems quite plausible that the probability that the world is good to some degree i when Ann is saved is just the probability that the world is good to that same degree when Bob is. Furthermore, for very high levels of goodness, it seems plausible that it is more probable that the world is good to that degree when both are saved than when just one is. Hence, it seems the desirabilities that are yielded by application of DAE conform sufficiently well to those furnished by intuition so as to disarm the counterexample.

Moreover, DAE can be seen to be nothing more than a reformulation of Jeffrey's desirability equation, given the existence of the goodness-level partition. For it follows from Jeffrey's formula that $V(A) = \sum_i V(A \wedge G_i)P(G_i | A)$, since the G_i form a partition of the space of possible worlds. But:

$$\begin{aligned} V(A \wedge G_i) &= \sum_{w_j \in A \wedge G_i} V(w_j)P(w_j | A \wedge G_i) \\ &= \sum_{w_j \in A \wedge G_i} i.P(w_j | A \wedge G_i) \end{aligned}$$

in virtue of the fact that by definition, $V(w_j) = i$ for all $w_j \in G_i$. So since $\sum_{w_j} P(w_j | A \wedge G_i) = 1$, it follows that $V(A \wedge G_i) = i$. Hence, Jeffrey's formula entails DAE, given the existence of the G_i -partition.

But now we run into a very interesting problem. We have argued that DAE is implied by Jeffrey's formulation of desirability (given the existence of the goodness-level partition). We have also seen that Invariance is false for desirability: The conditional desirability of A given A equals the desirability of the tautology and is typically not the same as the unconditional desirability of A . In contrast, DAE seems to entail Invariance: Learning that A does not change the value of $\sum_i i.P(G_i | A)$, since that would require that $P(G_i | A) \neq P_A(G_i | A)$. And this inequality can never rationally hold, if rational agents change their beliefs by Bayesian updating.¹³

Another way to put the problem, is that it seems that DAE cannot be maintained

¹³Here is a proof of that $P(\cdot | A)$ does not change when we conditionalise on A : If $P_A(B) = P(B | A) = P(A \wedge B) / P(A)$, then $P_A(B | A) = P_A(B \wedge A) / P_A(A) = [P(A \wedge B \wedge A) / P(A)] / P(A \wedge A) / P(A) = P(B \wedge A) / P(A) = P(B | A)$.

as an agent learns new propositions. Recall that the DAE equation states that $V(A) = \sum_i i.P(G_i | A)$. But we know that the left hand side of this equation normally changes as an agent learns the proposition A . But the same is not true for the right hand side of this equation (assuming that agents respond to learning by Bayesian updating). Hence, if the equation is satisfied before an agent learns that A , then it cannot still be satisfied after the agent learns this proposition.

In next section we offer a solution to this problem. Before doing so, we should point out that a causal version of Desire-as-Expectation is similarly entailed by causal decision theorists' concept of efficacy value. Recall that causal decision theorists in general say that the causal efficacy value of a proposition A is given by: $\sum_j V(w_j).P_A^\square(w_j)$, where P_A^\square is meant to be a variable that can represent whatever probability causal decision theorists take to be relevant when evaluating the choice worthiness of A (e.g. objective chance conditional on A , the image of P on A , etc.). But then we have:

$$\begin{aligned} U(A) &= \sum_j V(w_j).P_A^\square(w_j) \\ &= \sum_i i \sum_{w_j \in G_i} P_A^\square(w_j) \\ &= \sum_i i.P_A^\square(G_i) \end{aligned}$$

So causal efficacy value entails a causal Desire-as-Expectation thesis. This thesis, like the evidential one, is not undermined by our counterexample against the simple DAE, precisely because it allows for different degrees of goodness.

5.1 Conditioning with Indexical Propositions

The problem we face is the following. Jeffrey's theory implies both that the DAE thesis is true and that Invariance is false. But if agents revise their beliefs by Bayesian conditionalisation then DAE seems to imply Invariance. The aim of this section is to find a way out of this dilemma by giving a plausible explanation for why learning that A is the case changes the conditional probabilities for the G_i in such a way that the DAE equation can be sustained when updating on A .

Let's first get clear about what intuitively goes on when we change our views about some proposition. As before, let a proposition be a set of possible worlds. It is generally assumed that the content of a proposition—i.e., what worlds make it up—remains fixed when an agent changes her mind. For instance, when an agent changes her probability for some proposition A , it is assumed that the worlds making up A remain the same, while their probabilities change. Similarly, when an agent changes her mind about the desirability of A , this is understood as a change in the probability distribution over the worlds within A ; either from the less desirable worlds in A to the more desirable ones, or vice versa, and not as a change in the worlds constituting A .

In contrast, it *not* correct to assume that the contents of the goodness-level propositions—the G_i s—are fixed or invariant under changes in our evaluation of the desirability (and hence, by DAE, the goodness) of other propositions. Recall that the goodness-level propositions partition the space of possible worlds. When we change our view about the goodness of some proposition A , it gets a new place within this partition. Suppose for instance that A was originally a subset of G_j and that after we change our mind or get new information about its goodness it becomes a subset of G_k . Since the content of A has not changed—it is still made up of the same worlds as before—this means that the contents of G_k and G_j must have changed—the former contains worlds it didn't contain before, whereas the latter now does not contain worlds it did contain before.

To take an example, suppose we think that it would be very bad if the Liberal Democrats won the next UK general election; call this proposition L . However, having heard the leader of the party set out its policies, we change our mind, and conclude that the party isn't as bad as we thought. Now the content of the proposition L has not changed, although the probability distribution within L has shifted (compared to before, we are now more confident that one of the better worlds in L is actual if L is true). But crucially, the situation of L within the goodness partition has changed, and now occupies one of the 'better' regions in the goodness-level partition than before.

The upshot of this is that the G_i s are not strictly propositions, qua sets of possible worlds, but functions (of desirabilities) taking propositions as values. And so when an agent learns that A is true she must revise three things: Her probabilities for the possible

worlds, her desirabilities for these worlds, and the contents of the proposition-valued functions G_i . In particular, when she learns that A she not only revises (upwards) the probability of any world consistent with A but also revises (in the direction of the value of the tautology) its desirability. As a result of doing so the worlds consistent with A will come to belong to different goodness-level ‘propositions’ than before, which in turn will imply shifts in the conditional probabilities of (some or all) G_i s given A . Similarly, one should expect that for the G_i with a high (low) value of i , the conditional probability of G_i given A increases (decreases) when a person becomes more confident that A is good.

To spell this out formally, recall that P_A and V_A are the agent’s new probability and desirability functions after learning that A is true. Let the G_i be proposition-valued functions of desirability with $G_i(V) = \{w : V(w) = i\}$ and $G_i(V_A) = \{w : V_A(w) = i\}$. Note that $G_i(V_A) \neq G_i(V)$ when Invariance fails. That is, although both express the fact that the world has goodness of level i , the worlds making them true are different. (We could say that the sentences expressing $G_i(V)$ and $G_i(V_A)$ have the same intensional content but different extensional contents.) Now as conditionalisation on A does not change any conditional probabilities given A , Bayesians require that $P_A(G_i(V) | A) = P(G_i(V) | A)$ and $P_A(G_i(V_A) | A) = P(G_i(V_A) | A)$. But because of the (possible) difference in content between $G_i(V_A)$ and $G_i(V)$, we have: $P(G_i(V_A) | A) \neq P(G_i(V) | A)$.

This shows that, strictly speaking, if we want to make a version of the Desire-as-Expectation thesis compatible with the failure of Invariance, we should replace Broome’s thesis with:

Thesis 4 (Desire-as-Expectation*). *For any A and according to any rational agent:*

$$V(A) = \sum_i i.P(G_i(V) | A) \tag{9}$$

Then, contrary to what seemed to be the case, Bayesian conditioning is perfectly consistent with the fact that learning may change the conditional probability of the world being good to some degree given A . Similarly, once the three-fold effect of learning that

A is true is recognised we see that our version of DAE implies that:

$$\begin{aligned}
 V_A(A) &= \sum_i i.P_A(G_i(V_A) | A) \\
 &= \sum_i i.P(G_i(V_A) | A) \\
 &\neq \sum_i i.P(G_i(V) | A) = V(A)
 \end{aligned}$$

So Desire-as-Expectation* and Bayesian conditioning are jointly consistent with a denial of Invariance.

Our findings in this section provide support for a suggestion made by Alan Hájek and Philip Pettit (2004). They suggest that goodness is indexical in the same way we have said it must be—i.e., partly a function of a person’s attitudes—and they show that Lewis’ argument against DAB then loses its bite. Moreover, they explain why various meta-ethical views are committed to this indexicality. There are, nevertheless, important differences between our discussion of this issue and Hájek and Pettit’s. First, they accept Lewis’ argument against DAB as sound and suggest an indexical DAB thesis to avoid his result. We, on the other hand, have argued that Lewis’ argument is not sound, but that we nevertheless need an indexical account of goodness to save the Desire-as-Expectation thesis. Secondly, unlike Hájek and Pettit, we have shown that unless goodness is indexical, Jeffrey’s decision theory leads to contradiction, since it entails both the truth of Desire-as-Expectation and the falsity of Invariance, which is inconsistent unless goodness is indexical. Finally, the indexical thesis they suggest does not assume that there are multiple degrees of goodness, and is therefore refuted by the counterexample we discussed in section 4.

5.2 Anti-Humeans and subjectivists can also accept DAE

We have seen that our version of Desire-as-Expectation is not only consistent with Jeffrey’s decision theory, but an implication of it. But this might give rise to the worry that, contrary to Broome’s claim, Desire-as-Expectation is not properly anti-Humean, since on any defensible version of DAE, the goodness level partition always changes when the agent’s desires change.

This worry is, however, unnecessary. Anti-Humeans need not say that people are motivated directly through their beliefs about the good without these beliefs affecting their desires. (If they did say that, then they would have to deny that the structure of desire is captured by Jeffrey's formula.) In fact, the anti-Humean view can be characterised as precisely the idea that because rational people's beliefs about the good determine their desires, these beliefs determine what people are motivated to do. Broome, for instance, characterises the anti-Humean view thus:¹⁴

Sometimes, we do what will serve the good according to our beliefs about what would be good together with our other beliefs—no desire, *other than desires which result from beliefs alone*, need enter into it" (Broome (1991): 266, emphasis added)

If this is how we understand the anti-Humean view, then it is not a problem that a person's expectation of the good changes with her desires, if this change in desires is brought about by a change in beliefs. And that is, for instance, exactly what happens when expectation of good changes because a desire for a proposition has changed as a result of a change in the proposition's probability. More generally, there is no reason why an anti-Humean could not endorse our idea that the contents of the goodness propositions change with an agent's desires, but argue that what grounds such a change is (often, at least) a change in the agent's normative beliefs.

Moreover, since we can interpret the $G_i(V)$ functions however we like, the Desire-as-Expectation thesis should be acceptable to subjectivists as well as objectivists about value. We could, for instance, interpret $G_i(V)$ as expressing the fact that my desires are satisfied to degree i . Then DAE states that we should desire a proposition to the degree that we expect our desires to be satisfied when the proposition is true. This is something that people might want to accept irrespective of where they belong in the Humean/anti-Humean and objectivist/subjectivist divide.

¹⁴Broome is in the quoted passage rephrasing Lewis' (1988: 324) characterisation of anti-Humean. According to Lewis, anti-Humeans say that the only desires that motivate people to act are those that are *identical* with beliefs.

6 A generalisation of Desire-as-Belief?

Instead of avoiding the counterexample we raised against the simple Desire-as-Belief thesis by switching to some version of Desire-as-Expectation, a defender of DAB might respond by generalising the original DAB thesis to a thesis that allows for multiple degrees of goodness. As we saw at the start of last section, our counterexample would not undermine such a generalisation of DAB. Moreover, anti-Humeans should, independently of any discussion of DAB, be skeptical of the idea that goodness comes only in two degrees. Lewis himself considered such a generalised version of DAB and made a similar argument against it to that which he made against the simple DAB thesis (Lewis (1988): 330-331). But this argument again assumed Invariance. Therefore, a version of DAB based on multiple levels of goodness is neither refuted by Lewis' argument nor by our counterexample. But is it consistent with the truth of Desire-as-Expectation*? We conclude the paper with an argument for their compatibility.

To state the general version of the Desire-as-Belief thesis more precisely, let \mathring{A}_i be the proposition that A is good to degree i .¹⁵ Then the thesis under consideration states that:

Thesis 5 (Generalised Desire-as-Belief (GDAB)). *For any A and according to any rational agent:*

$$V(A) = \sum_i i.P(\mathring{A}_i) \quad (10)$$

If (our version of) the DAE thesis is true, then GDAB entails that:

$$\sum_i i.P(\mathring{A}_i) = \sum_i i.P(G_i(V) | A)$$

And that might seem very plausible. Indeed if the probability that A would be good to degree i equals the conditional probability that the world is good to degree i given that A is true—i.e. if $P(\mathring{A}_i) = P(G_i(V) | A)$ —then this implication must hold.

But now we might seem to be faced with another negative result due to David Lewis: his famous ‘triviality result’ against the so-called Adams’ thesis (see e.g. Lewis (1976),

¹⁵A generalised causal version of DAB could also be considered. Everything we say here about the evidential GDAB also holds for a causal GDAB, since, as we saw in last section, causal decision theory entails a causal DAE thesis.

Lewis (1986), and Hájek and Hall (1994)). Lewis' main target in his triviality argument was the idea that the probability of an indicative conditional, $A \rightarrow B$, is identical to the conditional probability of B given A . However, his result is easily generalised to a refutation of any claim of the form that for any probability function P and propositions A and B , there exists a proposition C such that $P(A | B) = P(C)$. So in particular it refutes the claim that there exists a proposition \mathring{A}_i such that $P(\mathring{A}_i) = P(G_i(V) | A)$. For this reason, Broome (1991) insists that we must resist the temptation to identify the probability that A is good to some degree i with the conditional probability that the world is good to degree i given A .

However, the aforementioned triviality results depend on taking the contents of the propositions in question to be fixed. In particular, learning some proposition is not supposed to change the content of any other proposition (nor of that same one). But in the last section we argued that, first, DAE entails that this is precisely what happens with the goodness-level 'propositions' when one changes ones mind about how good some proposition is, and, second, that this is what must happen if DAE is not to imply the false Invariance principle.

The contents of the \mathring{A}_i propositions must not be fixed either, for similar reasons. These propositions also partition the space of possible worlds, but when a person learns that A is true, for instance, then unless A was already considered neither more nor less desirable than the tautology, A 's place within the \mathring{A}_i partition must (if GDAB is true) shift such that its expected goodness becomes equal to that of the tautology (recall our discussion from section 3.2.). But then since the content of A does not change when an agent learns that the proposition is true, the content of some \mathring{A}_i proposition must change when an agent learns this. So the contents of both \mathring{A}_i and $G_i(V_A)$ change as an agent learns new information. The upshot is that the triviality arguments like those that have been taken to undermine Adams' thesis do not invalidate the above argument for GDAB.

A similar treatment can be given of another apparent problem for the generalised DAB thesis, namely that it seems that, given DAE, the value propositions \mathring{A}_i must be probabilistically independent of A , which then again entails the false Invariance

principle. For since $\sum_i i.P(G_i(V_A) \mid A)$ might seem invariant under changes in the probability of A , it follows that if both DAE and GDAB are true, then $\sum_i P(\mathring{A}_i)$ must also be invariant under changes in A . But then by GDAB, $V(A)$ must also be invariant under changes in the probability of A , which we have seen to be inconsistent with Jeffrey's understanding of desirability.

The above worry is mistaken however. In last section we showed that $\sum_i i.P(G_i(V_A) \mid A)$ is *not* invariant under changes in the probability of A . It only seems to be so because of an implicit assumption that the contents of the goodness-level 'propositions' are fixed. But we have shown, first, that one should not make that assumption, and, second, that our formulation of DAE does not entail it. But then if both DAE and GDAB are true, $\sum_i P(\mathring{A}_i)$ may also change when an agent learns that A is true, just as the falsity of Invariance entails. In sum, the Desire-as-Expectation* thesis does not contradict the Generalised Desire-as-Belief thesis, given the indexical nature of the goodness-level propositions.

7 Concluding remarks

Contrary to what David Lewis thought, Bayesian decision theory does not rule out the possibility that there is a fixed quantitative relationship between a rational person's desires and her evaluative beliefs. Although the simple version of the Desire-as-Belief thesis that he proposed is refuted by our counterexample, the more plausible Desire-as-Expectation* thesis—which should be acceptable to Humeans, anti- Humeans, subjectivists and objectivists—is not only consistent with Bayesian decision theory but entailed by it. And this thesis, in turn, is consistent with a version of DAB that allows for multiple goodness levels. Nor, contrary to appearances, do these theses imply the dubious Invariance principle since propositions expressing goodness claims must themselves have contents that vary. So Bayesian decision theory, it would seem, is general enough to allow for a range of different theories of value and human motivation.

References

- Bradley, R. and C. List (2009). Desire-as-belief revisited. *Analysis* 69(1), 31–37.
- Broome, J. (1991). Desire, belief and expectation. *Mind* 100(2), 265–267.
- Byrne, A. and A. Hájek (1997). David Hume, David Lewis, and decision theory. *Mind* 106(423), 411–428.
- Collins, J. (1988). Belief, desire, and revision. *Mind* 97(387), 333–342.
- Costa, H. A., J. Collins, and I. Levi (1995). Desire-as-belief implies opinionation or indifference. *Analysis* 55(1), 2–5.
- Daskal, S. (2010). Absolute value as belief. *Philosophical Studies* 148(2), 221–229.
- Hájek, A. and N. Hall (1994). The hypothesis of the conditional construal of conditional probability. In *Probability and Conditionals: Belief Revision and Rational Decision*. Cambridge University Press.
- Hájek, A. and P. Pettit (2004). Desire beyond belief. *Australasian Journal of Philosophy* 82(1), 77–92.
- Jeffrey, R. (1990/1983). *The Logic of Decision*. The University of Chicago Press (paperback edition).
- Lewis, D. (1976). Probabilities of conditionals and conditional probabilities. *Philosophical Review* 85(3), 297–315.
- Lewis, D. (1981). Causal decision theory. *Australasian Journal of Philosophy* 59(1), 5–30.
- Lewis, D. (1986). Probabilities of conditionals and conditional probabilities II. *Philosophical Review* 95(4), 581–589.
- Lewis, D. (1988). Desire as belief. *Mind* 97(387), 323–32.
- Lewis, D. (1996). Desire as belief II. *Mind* 105(418), 303–313.
- Nissan-Rozen, I. (2015). A triviality result for the “Desire by Necessity” thesis. *Synthese* (forthcoming).

- Nozick, R. (1969). Newcomb's problem and two principles of choice. In N. Rescher (Ed.), *Essays in Honor of Carl G. Hempel*. Reidel.
- Oddie, G. (1994). Harmony, purity, truth. *Mind* 103(412), 451–472.
- Oddie, G. (2001). Hume, the BAD paradox, and value realism. *Philo* 4(2), 109–122.
- Plato (1956/360 BC). *Symposium*. Pearson (translated by Benjamin Jowett).
- Russell, B. (1905). On denoting. *Mind* 14(56), 479–493.
- Smith, M. (1987). The Humean theory of motivation. *Mind* 96(381), 36–61.
- Smith, M. (1994). *The Moral Problem*. Wiley-Blackwel.
- Weintraub, R. (2007). Desire as belief, Lewis notwithstanding. *Analysis* 67(2), 116–122.
- Williams, J. R. G. (2010). Counterfactual desire as belief. Unpublished manuscript.