

Konek & Levinstein: The Foundations of Epistemic Decision Theory

KEVIN DORST

11.9.15

I. (Alethic) Epistemic Decision Theory

Accuracy as *the* fundamental epistemic good: “The higher your credence in truths and the lower your credence in falsehoods, the better off you are all epistemic things considered... [So] norms of epistemic rationality, on this view, have their binding force in virtue of the following fact: they are good means toward the end of securing accuracy.” (1)

Make “good means toward the end” precise by co-opting resources from practical decision theory. But (Greaves thinks) the Imps result will then be inevitable, since a reasonable accuracy metric will assess accuracy globally (via weightings), which means it will be open to tradeoffs amongst propositions (2).

They Claim: IMPS requires re-thinking of foundations of EpDT, but doesn’t show epistemic rationality to be non-consequentialist. In particular, there are **two ways of evaluating how accurate Emily is:**

- 1) Evaluate her epistemic *state*: how close her credences are (at t) to the truth-values (at t).
- 2) Evaluate her epistemic *act*: how much accuracy her coming-to-occupy state c (denoted \bar{c}) *produced*.

But epistemic rationality tracks (2), not (1), since “epistemic states, rather than acts, have the epistemically interesting direction of fit, viz., mind-to-world.” (3)

∴ EpDT still sanctions the epistemic *act* of taking the bribe, but the *state* is the locus of evaluation when it comes to rationality.

Q: Can the direction-of-fit metaphor sustain the distinction? Here’s Emily; she’s wondering what to believe; she’s in a decision-theoretic frame of mind (“I’m out to get accuracy!”); she knows that she’ll get more of it if she takes the bribe. What, exactly, are we supposed to say to convince her not to take it?

II. Praxic and Epistemic Theories of the Good

A. General theory of rational preferences:

Agents ought to line up their preferences over options (acts, epistemic states) with their *unconditional best estimates* of the value (practical, epistemic) of those options.

Q: This seems to be a straightforward endorsement of Accurate Alice. Is she irrational since she’s not following (say) the Principal Principle? She’s certain the chance is .5 that p (and it is), and she’s certain that p (because p in fact will happen). Yet she’s (correctly) sure that conforming to PP will make her less accurate; so if the normative force of such principles stems from their ability to promote accuracy, what has EpDT got to say against her?

They say (1) and (2) clearly come apart in this case: take an agent certain in $\neg C_0 \wedge C_1 \wedge \dots \wedge C_{10}$; “she knows the credence function that assigns 1 to C_0 and C_1, \dots, C_{10} is more accurate than her own. However *were she to adopt that state*, she would end up less accurate than she currently is.” (2)

Q: Is that right? She has credence 0 in C_0 ; how could she know it’s true?

Consequentialist Evidentialism, I think, gives more distance. The theory is thoroughly deontological — you’re (epistemic) job is *not* to go get accuracy, so there’s no reason for bribes to move you.

This will in fact be consistent with EDT and CDT; it’s all about getting subtle with the values.

Direction-of-fit metaphor; Anscombe's shopper/detective story.

Praxic Good: An action A is prudentially valuable at world w , relative to a state of desire D , to the extent that A makes w satisfy D , by causally influencing it in the right way.

Epistemic Good: A doxastic state B is epistemically valuable at a world w to the extent that B is close to the truth at w .

How do we make these precise?

Epistemic case is simple: use unconditional estimate of accuracy, à la Savage. This is because (basically) we don't care about how your epistemic states *affect* the world; we just care about how close you expect them to get to it.

Practical case involves the notions of indicative ($c(w|X)$) vs. subjunctive ($c(w||X)$) supposition. The first measures how well X correlates with w . The second measures how well X causally contributes to w .

How to measure causal impact? "Imaged Bayes factor": $\frac{c(w|X)}{c(w)}$. The idea is then to use this factor to weight utility u of any given world w . Thus the value of a world w , relative to an act A , depends on (i) how desirable w is, and (ii) how well A does at bringing it about that w (5). More precisely, we'll have

$$V_A(w) = \frac{c(w|A)}{c(w)} \cdot u(w)$$

B. Comparison to Joyce

Savage says prefer based on your unconditional estimated value for various acts, where value is just $u(w)$.

Joyce says Savage evaluates from the *wrong epistemic perspective* — acts ought to be evaluated on the supposition that they are performed (7). General Joycean theory of preference: Prefer A to B iff:

$$\sum_i (c(S_i||A)u(o[A, S_i])) \geq \sum_i (c(S_i||B)u(o[B, S_i]))$$

K&L disagree: we shouldn't characterize the difference between EDT and CDT in terms of their epistemic perspectives; instead they have *different theories of (praxic) good*. "They disagree, in the first instance, about which quantity to estimate, for the purpose of evaluating actions..." (9). Thus we get:

$$\text{CDT theory of value: } V_A(w) = \frac{c(w|A)}{c(w)} \cdot u(w)$$

$$\text{EDT theory of value: } V_A(w) = \frac{c(w|A)}{c(w)} \cdot u(w)$$

Since credal states are mind-to-world, they *aren't* evaluated in this "contribute to good outcomes" way, so the utilities shouldn't be weighted

You can influence the world in a variety of ways. "But — and this is the crucial point — influencing the world in (epistemically) good or bad ways is not what makes [beliefs] epistemically valuable. What makes them epistemically valuable — the primary source of all epistemic-things-considered value — is just accuracy." (4)

Example of waking up on the roof w/ a rifle and a despot in view. Indicative supposition: less likely despot will die; subjunctive supposition: more likely he'll die.

This, in generally, requires act/state independence.

EDT then says that $c(S_i|||A) = c(S_i|A)$, while CDT says $c(S_i|||A) = c(S_i|A)$.

The basic thought here is that when we're evaluating A , we weight the utility of w by how well A contributes to w obtaining; CDT understands "contributes" causally; EDT understands "contributes" evidentially.

by the “acts” contribution to making it the case that w .

“Greaves’ concerns about accuracy-first epistemology result from running these very different sorts of evaluations — evaluations of epistemic states and acts — together. Carefully separating them out is the key to seeing that accuracy-first epistemology does not sanction epistemic bribe-taking.” (11)

III. Solving the puzzles

They (crucially, it seems) use two principles to solve the puzzles:

Principal Principle (PP)

Deference to Chance (DtC): if an agent has credences that estimate the inaccuracy of p to be x , and she is certain that the chance function estimates its inaccuracy to be y ($y \neq x$, so she violates PP), then she calculates estimates of inaccuracy relative to the chances. (11)

Now, we always bet using *estimates*. Usually those will be expectations from credences, but in cases described in DtC they won’t.

Q: How can an accuracy-firster help themselves to such principles? They cite Pettigrew-style justifications of PP (via chance-dominance) (11-12), but *given the decision-theoretic interpretation*, why should such an agent defer to the chances?

Q: Also, I start to get confused at how the pieces are working here. It looks like DtC only kicks in when PP fails; so given that they have PP, why do they need it?

A. *Imps*

Let Em_x be Emily’s doppelgänger w/ credence x in C_0 , that there’s a child in front of her.

Em_1 ought to prefer her own credence because (I think?) (i) it satisfies PP, (ii) we have a strictly proper scoring rule, and (iii) we’re using unconditional estimates of accuracy.

If $x \neq 1$, Em_x will prefer to prefer the *chances* to her actual credence, since (i) she knows $ch(C_0) = 1$, so (ii) by DtC she’ll use the chances to calculate her estimates, and (iii) by strict propriety the chances will prefer themselves.

In particular, then, given that she *has* credence 1, looking at any other credence will still have the $\frac{1}{2}$ uncertainty in each C_i . The estimate ignores the causal impact that changing your credence will have.

“Less formally, since Emily ought to prefer to be in whatever credal state is, in her best estimate, most valuable, i.e. most accurate; and since she treats chance’s best estimates as her own (by DtC); and moreover since she is certain that the true chance function ch_E estimates itself to be most accurate, she ought to prefer ch_E — a credal state which assigns 1 to the propositions that there’s a child before her — to any other credal state b , including her own c_x ” (12-13).

Q: Okay, how does this work?

(i) How much work is DtC doing for the theory? DtC (given propriety) guarantees that wherever she starts, she'll go to $c(C_0) = 1$. But if so, why did we need the first half of the paper? *Answer(?)*: Without re-thinking the theory of preference, even with propriety the chances would not prefer themselves. Propriety is a constraint using *unconditional* estimates; but if we use estimates conditional on some act then *even if her credences match the chances*, they'll still sanction the epistemic bribe.

(ii) Let's think hard about this DtC stuff, because it really is crucial to the solution. Is there good enough reason to accept it, from an accuracy-firster point of view? You're certain in p . You also know the objective chances of p are .5. But why should you care? After all, if you stick with you're credence you'll definitely get it right, whereas if you conform to the chances you'll get it wrong. It might help to distinguish two ways the world could be indeterministic.

(1) Aristotelian indeterminism: there is no fact of the matter, now, whether they'll be a sea battle tomorrow.

(2) Quantum indeterminism: the objective chance of there being a sea battle tomorrow is .5.

Now *if* objective chances take the form of (1), then I could see why you need to defer: in some sense you *can't* outsmart the chances. But if it's just (2) (as it usually is, I think), why would you defer?

Put it this way: why isn't the fact that you're certain of p enough to show that you have "inadmissible information" (à la time travelers), and so nullify PP and DtC?

IV. What EpDT Recommend

If we have time: in evaluating *acts*, EpDT recommends shifting to credence 0 in C_0 . But in terms of *states* it recommends credence 1. Basically, they think that the direction-of-fit discussion shows that epistemic rationality is concerned with evaluating states.

The only sense in which there's a rational dilemma is the normal dilemmas between practical and epistemic rationality.

You *could* be a really weird agent, who only cared about accuracy. Then decision theory (epistemic or otherwise) would recommend taking epistemic bribes. But that's something everyone who likes practical decision theory has to accept.