

Against Tournament Decision Theory

Brian Weatherson

2020-11-02

Abelard Podgorski (2020) argues in favour of what he calls Tournament Decision Theory. The core idea is that rather than assigning each option a value, and then choosing the option with the greatest value, we devise a formula for making binary comparisons between options. Podgorski notes there are a few things that you can do with these binary comparisons, but a minimal commitment of his theory is that if an option is better than each alternative by the lights of the preferred pairwise comparison, then it does best. We'll show that this leads to implausible results - in particular it leads to choosing dominated options.

To set out his theory we need one bit of familiar conceptual machinery. Let $V_A(B)$ be what value the decider expects B would have were they to choose A . So in a case where we expect a demon to predict our decision, this is the value in the cell of the table where B is chosen and A is predicted. Then say that X is pairwise better than Y in case:

$$V_X(X) + V_Y(X) > V_X(Y) + V_Y(Y)$$

And, as noted in the first paragraph, a sufficient condition for an option being best is that it does better than any other option in any pairwise comparison. That's enough to get a counterexample to the theory.

Our agent has three choices, and if they choose any one of them with certainty, they are very confident that a demon will correctly predict their choice. And the payoffs for each choice/prediction pair are given by the following table.

	A predicted	B predicted	C predicted
A chosen	0	3	5
B chosen	2	2	2
C chosen	5	3	0

The first thing to say about this example is that B is the worst thing to choose. Here are two reasons to think it is wrong, or at least what looks like two reasons.

First, it is dominated by the mixed strategy of playing A with probability 0.5 and C with probability 0.5. And one should never play dominated strategies, even if what does the dominating is a mixed strategy.

Second, there is no credence distribution over the Demon's predictions that makes B the choice with maximal expected utility. In the terminology of game theorists, B is not a Best Response - it doesn't maximise

expected utility no matter what one's beliefs. Agents should have probabilistically coherent credences, and maximise expected utility given their beliefs, and no agent that satisfies both conditions chooses *B*.

Now I've said these are two reasons to not choose *B*, and they appeal to different conceptual tools - the first appeals to anti-dominance reasoning, the second to utility maximisation. But there is a sense in which they are the same reason. As Pearce (1984) shows, a strategy is undominated (in the sense I'm using here) just in case there is some credence distribution that makes it a utility maximising choice. It would take us too far afield to work through whether this means the 'two' reasons from the last two paragraphs are really the same reason. Despite their mathematical equivalence, I suspect they are convincing to different theorists, so I've included both. And they point to the same conclusion: *B* is a uniquely bad choice.

The second thing to say about this example is that Podgorski's theory says that *B* is the unique rational choice. This can be seen simply applying the above formula to confirm that the theory judges *B* to be better than *A* and better than *C*. But it is useful to abstract to see just what has gone wrong here. Replace the values in the above table with variables as follows:

	A predicted	B predicted	C predicted
A chosen	v_{11}	v_{12}	v_{13}
B chosen	v_{21}	v_{22}	v_{23}
C chosen	v_{31}	v_{32}	v_{33}

The condition for *B* beating *A* in the pairwise comparison is that $v_{21} + v_{22} > v_{11} + v_{12}$. The condition for *B* beating *C* in the pairwise comparison is that $v_{22} + v_{23} > v_{32} + v_{33}$. And that's satisfied in our example in both cases, since $2 + 2 > 0 + 3$. But the more important thing to note is what's not there. The sufficient condition Podgorski gives does not even mention v_{13} or v_{31} . It is very implausible for anyone except a dedicated Evidential Decision Theorist to think that the values of v_{13} or v_{31} are irrelevant to a sufficient condition for choosing *B*.

The general lesson here is I think a quite general one. As William Harper (1986) showed, most of the problems that decision theorists worry about are basically games where the Demon's utility function is left implicit. The Demon wants to make a correct prediction, so we can turn an example into a game by just saying that the Demon's payoff is 1 if the prediction is correct, and 0 otherwise. Here's how the example I described above looks if we do this.

	A predicted	B predicted	C predicted
A chosen	0, 1	3, 0	5, 0
B chosen	2, 0	2, 1	2, 0
C chosen	5, 0	3, 0	0, 1

Now it is a matter of some dispute among game theorists just what either player should do in this game.

But one very plausible minimal condition is that each player should choose a strategy that is rationalisable in the sense of Bernheim (1984) and Pearce (1984).¹ Since any (reasonable) game has Nash equilibria, and all Nash equilibria are n-tuples of rationalisable strategies, this is guaranteed to not rule out all strategies.² And it's hard to see why it could be ever be good to choose a non-rationalisable strategy, since it is hard to see what coherent mental state could issue in such a choice.

There is an even weaker constraint that we could impose: only choose Best Responses. That is, we could require that choosers have probabilistically coherent credences over the possible states of the world, and that they maximise expected utility given those credences. Since all rationalisable strategies are Best Responses, but not vice versa, this is a strictly weaker requirement. But choosing *B* violates even this requirement.

Saying that one should only choose Best Responses, or even that one should only choose rationalisable strategies, isn't the end of the story. Nothing I've said here suggests an answer to the question of whether there is a further constraint beyond the constraint that one choose rationalisable strategies. And if there are such constraints, it doesn't suggest anything about what they should be. It is easy enough to see what some constraints could be. For example, we get interesting variants of CDT by requiring that players adopt strategies that are parts of Nash equilibria, or are evolutionarily stable. These variants resemble, though they aren't quite identical to, the view that Podgorski calls Deliberational CDT. Given the huge range of solution concepts that have been developed by game theorists over the past 40 years, Harper's idea of thinking of decisions as games suggests rather a lot of possible decision theories. But all of them should start from the basic idea that only rationalisable strategies are permissible, and that rules out Tournament Decision Theory.

Bernheim, B. Douglas. 1984. "Rationalizable Strategic Behavior." *Econometrica* 52 (4): 1007–28. <https://doi.org/10.2307/1911196>.

Harper, William. 1986. "Mixed Strategies and Ratifiability in Causal Decision Theory." *Erkenntnis* 24 (1): 25–36. <https://doi.org/10.1007/BF00183199>.

Pearce, David G. 1984. "Rationalizable Strategic Behavior and the Problem of Perfection." *Econometrica* 52 (4): 1029–50. <https://doi.org/10.2307/1911197>.

Podgorski, Aberlard. 2020. "Tournament Decision Theory." *Noûs* tbc (tbc): xx–xx. <https://doi.org/10.1111/nous.12353>.

¹A strategy is rationalisable if it survives the following process of iterated deletion. At each stage delete a strategy unless there is some probability distribution over the strategies of other players such that the expected utility of the strategy is maximal given that probability distribution.

²This is a respect in which this requirement, which looks a lot like the admonition to choose ratifiable strategies, is in fact much more plausible than a ratifiability requirement.