

Prisoners' Dilemma

PARADOX AND INFINITY

Benjamin Brast-McKie

April 29, 2024

Two Prisoners

Setup: Two separated prisoners are each offered \$1,000. They will be given an additional \$1,000,000 *iff* the other prisoner does not take the \$1,000.

- The prisoners' choices are causally independent.
- $P(\text{Take}_A \square \rightarrow \text{Take}_B) = P(\neg \text{Take}_A \square \rightarrow \text{Take}_B) = P(\text{Take}_B)$.
- $P(\text{Take}_A \square \rightarrow \neg \text{Take}_B) = P(\neg \text{Take}_A \square \rightarrow \neg \text{Take}_B) = P(\neg \text{Take}_B)$.
- We know that $P(\neg \text{Take}_B) = 1 - P(\text{Take}_B)$, but don't know $P(\text{Take}_B)$.
- Something similar may be said swapping 'A' and 'B' above.
- The prisoner's know everything except for the other's choice.
- What is it rational for prisoner A (similarly B) to do?

Dominant: Taking the \$1,000 is a *dominant strategy* for prisoner A (similarly B).

- Whether Take_B or not, $v(\text{Take}_A) > v(\neg \text{Take}_A)$ for prisoner A.
- We get the following alternatives:

	Take_B	$\neg \text{Take}_B$
Take_A	$(A, B : \$1,000)$	$(A : \$1,001,000), (B : \$0)$
$\neg \text{Take}_A$	$(A : \$0), (B : \$1,001,000)$	$(A, B : \$1,000,000)$

- The setup assumes that neither prisoner cares about the other.
- If the prisoners cared about each other, that would be a different case.

Predictor: Given the circumstances, each prisoner is a good predictor of the other.

- Take_A predicts that Take_B , i.e., $P(\text{Take}_B \mid \text{Take}_A)$ is high.
- Thus $P(\neg \text{Rich}_A \mid \text{Take}_A)$ is high since $\text{Take}_B \text{ iff } \neg \text{Rich}_A$.
- So if Take_A , then prisoner A has good reason to bet $\neg \text{Rich}_A$.
- We don't know what the probabilities $P(\text{Take}_A)$ or $P(\text{Take}_B)$.
- Does $\neg \text{Take}_A$ change the probability $P(\neg \text{Take}_B) = P(\text{Rich}_A)$?

Newcomb: Rich_A *iff* it is predicted that $\neg \text{Take}_A$ (by $\neg \text{Take}_B$).

- $\neg \text{Take}_B$ is a *prediction instance* (a way of predicting $\neg \text{Take}_A$).
- The predication amounts to probabilistic dependence (not causal).
- When the prediction happens does not matter to the case.
- Is the prisoners' dilemma a Newcomp problem?

Dominance Calculations

Expected Causal Utility: Recall that: (a) $\text{Rich}_A \text{ iff } \neg \text{Take}_B$; and (b) $\text{Rich}_B \text{ iff } \neg \text{Take}_A$.

- What are the *expected causal utilities* of Take_A and $\neg \text{Take}_A$?
- $ECU(\neg \text{Take}_A) = \$1,000,000 \times P(\neg \text{Take}_A \sqcap \rightarrow \text{Rich}_A) + \$0 \times P(\neg \text{Take}_A \sqcap \rightarrow \neg \text{Rich}_A)$
 $= \$1,000,000 \times P(\neg \text{Take}_B)$ by (a).
- $ECU(\text{Take}_A) = \$1,001,000 \times P(\text{Take}_A \sqcap \rightarrow \text{Rich}_A) + \$1,000 \times P(\text{Take}_A \sqcap \rightarrow \neg \text{Rich}_A)$.
 $= \$1,001,000 \times P(\neg \text{Take}_B) + \$1,000 \times P(\text{Take}_B)$ by (a).
 $= \$1,001,000 \times P(\neg \text{Take}_B) + \$1,000 \times (1 - P(\neg \text{Take}_B))$.
 $= \$1,000,000 \times P(\neg \text{Take}_B) + \$1,000$.
 $= ECU(\neg \text{Take}_A) + \$1,000$.
- Taking the money is better for prisoner A (and similarly for B).

Accuracy

Clash: The predication does not have to be very accurate for the expected utility calculation to clash with causal expected utility (i.e. $> .5005$).

- Suppose $P(\text{Take}_B | \text{Take}_A) = P(\neg \text{Take}_B | \neg \text{Take}_A) = .5006$.
- So $P(\text{Rich}_A | \neg \text{Take}_A) = P(\neg \text{Take}_B | \neg \text{Take}_A) = .5006$.
- And $P(\text{Rich}_A | \text{Take}_A) = P(\neg \text{Take}_B | \text{Take}_A) = 1 - P(\text{Take}_B | \text{Take}_A) = .4994$.
- $EV(\neg \text{Take}_A) = \$1,000,000 \times P(\text{Rich}_A | \neg \text{Take}_A) + \$0 \times P(\neg \text{Rich}_A | \neg \text{Take}_A)$
 $= \$1,000,000 \times P(\text{Rich}_A | \neg \text{Take}_A)$
 $= \$500,600$.
- $EV(\text{Take}_A) = \$1,001,000 \times P(\text{Rich}_A | \text{Take}_A) + \$1,000 \times P(\neg \text{Rich}_A | \text{Take}_A)$
 $= \$1,000,000 \times P(\text{Rich}_A | \text{Take}_A) + \$1,000$
 $= \$500,400$.
- Even if prisoner A is an inaccurate predictor of prisoner B , the expected utility and expected causal utility calculations are bound to come apart.

Upshot

Common: Newcomb's problem is fanciful, but prisoners' dilemmas are common.

- Prisoners' dilemmas support *causal decision theory* on their own.
- No need to appeal to Newcomb cases to motivate CDT.

Comparison: Should a oneboxer also avoid taking the money?

- Does comparing the cases put any pressure on the oneboxer to twobox?