Epistemic Game Theory Lecture 6

Eric Pacuit

University of Maryland, College Park pacuit.org epacuit@umd.edu

March 10, 2014

Newcomb's Paradox

Two boxes in front of you, A and B.

Box A contains \$1,000 and box B contains either \$1,000,000 or nothing.

Newcomb's Paradox

Two boxes in front of you, A and B.

Box A contains \$1,000 and box B contains either \$1,000,000 or nothing.

Your choice: either open both boxes, or else just open *B*. (You can keep whatever is inside any box you open, but you may not keep what is inside a box you do not open).

Newcomb's Paradox



A very powerful being, who has been invariably accurate in his predictions about your behavior in the past, has already acted in the following way:

- 1. If he has predicted that you will open just box *B*, he has in addition put \$1,000,000 in box *B*
- 2. If he has predicted you will open both boxes, he has put nothing in box *B*.

What should you do?

R. Nozick. Newcomb's Problem and Two Principles of Choice. 1969.

Orthodox Bayesian: It is a problem of act-state dependence (1-box)

- Orthodox Bayesian: It is a problem of act-state dependence (1-box)
- Causal Decision Theory: expected utility involves probabilities of causal counterfactuals (2-box)

- Orthodox Bayesian: It is a problem of act-state dependence (1-box)
- Causal Decision Theory: expected utility involves probabilities of causal counterfactuals (2-box)
- No Acyclic Reasons: reasoning cannot refer to the act of choice in an essential way (2-box)...plus some "mental gymnastics" (1-box)

- Orthodox Bayesian: It is a problem of act-state dependence (1-box)
- Causal Decision Theory: expected utility involves probabilities of causal counterfactuals (2-box)
- No Acyclic Reasons: reasoning cannot refer to the act of choice in an essential way (2-box)...plus some "mental gymnastics" (1-box)
- "Tickle": Pr(page box contain \$0 | T & 1-box) =
 Pr(page box contain \$0 | T & 2-box) (2-box)

- Orthodox Bayesian: It is a problem of act-state dependence (1-box)
- Causal Decision Theory: expected utility involves probabilities of causal counterfactuals (2-box)
- No Acyclic Reasons: reasoning cannot refer to the act of choice in an essential way (2-box)...plus some "mental gymnastics" (1-box)
- "Tickle": Pr(page box contain \$0 | T & 1-box) =
 Pr(page box contain \$0 | T & 2-box) (2-box)
- Evidential Decision Theory: decisions to act provides evidence for the consequences (1-box)

- Orthodox Bayesian: It is a problem of act-state dependence (1-box)
- Causal Decision Theory: expected utility involves probabilities of causal counterfactuals (2-box)
- No Acyclic Reasons: reasoning cannot refer to the act of choice in an essential way (2-box)...plus some "mental gymnastics" (1-box)
- "Tickle": Pr(page box contain \$0 | T & 1-box) =
 Pr(page box contain \$0 | T & 2-box) (2-box)
- Evidential Decision Theory: decisions to act provides evidence for the consequences (1-box)
- Ratifiability: decision makers must assess the act in light of the decision to perform it and only choose acts that are self-ratifiable (1-box)

Causal Decision Theory

A. Egan. *Some Counterexamples to Causal Decision Theory*. Philosophical Review, 116(1), pgs. 93 - 114, 2007.

Smoking Lesion: Susan is debating whether or not to smoke. She knows that smoking is strongly correlated with lung cancer, but only because there is a common cause a condition that tends to cause both smoking and cancer. Once we fix the presence or absence of this condition, there is no additional correlation between smoking and cancer. Susan prefers smoking without cancer to not smoking without cancer, and prefers smoking with cancer to not smoking with cancer. Should Susan smoke? Is seems clear that she should. In The Smoking Lesion there is a strong correlation between smoking and getting cancer, despite the fact that smoking has no tendency to cause cancer, due to the fact that smoking and cancer have a common cause. In The Smoking Lesion there is a strong correlation between smoking and getting cancer, despite the fact that smoking has no tendency to cause cancer, due to the fact that smoking and cancer have a common cause. Still, since Susan's p(CANCER | SMOKE) is much higher than her p(CANCER | NOT SMOKE), EDT assigns not smoking a higher value than smoking. And this seems wrong. The Psychopath Button: Paul is debating whether to press the 'kill all psychopaths' button. It would, he thinks, be much better to live in a world with no psychopaths.

The Psychopath Button: Paul is debating whether to press the 'kill all psychopaths' button. It would, he thinks, be much better to live in a world with no psychopaths. Unfortunately, Paul is quite confident that only a psychopath would press such a button.

The Psychopath Button: Paul is debating whether to press the 'kill all psychopaths' button. It would, he thinks, be much better to live in a world with no psychopaths. Unfortunately, Paul is quite confident that only a psychopath would press such a button. Paul very strongly prefers living in a world with psychopaths to dying. Should Paul press the button? (Set aside your theoretical commitments and put yourself in Pauls situation. Would you press the button? Would you take yourself to be irrational for not doing so?)

 $p(\text{press button } \Box \rightarrow \text{dead}) = 0.001$

 $p(\text{press button } \square \rightarrow \text{live in a world without psychopaths}) = 0.999$

This is because Paul either is or is not a psychopath, and the probability of the two possibilities does not depend on what he decides to do.

Press Button: $p(\text{press button} \Box \rightarrow \text{dead}) \cdot u(\text{dead}) + p(\text{press button} \Box \rightarrow \text{live in a world without psychopaths}) \cdot u(\text{live in a world without psychopaths}) = (0.001 \cdot -100) + (0.99 \cdot 1) = 0.89$

Do Not Press Button: $p(\text{do not press button } \Box \rightarrow \text{live in a world with psychopaths}) \cdot u(\text{live in a world with psychopaths}) = 1 \cdot 0 = 0$

Death in Damascus

A man in Damascus knows that he has an appointment with Death at midnight. He will escape Death if he manages at midnight not to be at the place of his appointment. He can be in either Damascus or Aleppo at midnight.

Death in Damascus

A man in Damascus knows that he has an appointment with Death at midnight. He will escape Death if he manages at midnight not to be at the place of his appointment. He can be in either Damascus or Aleppo at midnight. As the man knows, Death is a good predictor of his whereabouts. If he stays in Damascus, he thereby has evidence that Death will look for him in Damascus. However, if he goes to Aleppo he thereby has evidence that Death will look for him in Aleppo.

Death in Damascus

A man in Damascus knows that he has an appointment with Death at midnight. He will escape Death if he manages at midnight not to be at the place of his appointment. He can be in either Damascus or Aleppo at midnight. As the man knows, Death is a good predictor of his whereabouts. If he stays in Damascus, he thereby has evidence that Death will look for him in Damascus. However, if he goes to Aleppo he thereby has evidence that Death will look for him in Aleppo. Wherever he decides to be at midnight, he has evidence that he would be better off at the other place. No decision is stable.

A. Gibbard and W. Harper. *Counterfactuals and Two Kinds of Expected Utility*. In Ifs: Conditionals, Belief, Decision, Chance, and Time, pp. 153190, 1978.

Ratifiability

The notion of ratifiability is applicable only where, during deliberation, the agent finds it conceivable that he will not manage to perform the act he finally decides to perform, but will find himself performing one of the other available acts instead...The option in question is ratifiable or not depending on whether or not the expected desirability of actually carrying out each of the alternatives (in spite of having chosen to carry out a different option, as hypothesized) (Jeffrey, 1983, pgs. 18-20)

The crucial distinction is between an act and a decision to perform the act.

Before performing an act, an agent may assess the act in light of a decision to perform it. Information the decision carries may affect the act's expected utility and its ranking with respect to other acts.

Decision makers should make self-ratifying, or ratifiable, decisions.

Two Forms of Ratificationism

- As an *elimination rule*: ratificationism requires you to reject all unratifiable acts, and to then choose among the ratifiable alternatives.
- As an equilibrium rule: ratificationism requires you to choose an act that is ratifiable relative to the beliefs and desires you will have when your deliberations cease ("reflective equilibrium").

Causal Rationality in Games

O. Board. *The Equivalence of Bayes and Causal Rationality in Games.* Theory and Decision, 61, pgs. 1 - 19, 2006.

Aumann Model

 $\langle W, \{R_i\}_{i\in N}, \{f_i\}_{i\in N}, \{p_i\}\rangle$

- W is a (finite) set of states.
- For each *i* ∈ *N*, *R_i* is a relation on *W R_i* is serial (and transitive, Euclidean, etc.)
 Let *R_i(w)* = {*v* | *w R_i v*}

► For each
$$i \in N$$
, $f_i : W \to S_i$
 $f_{-i}(w) = \langle f_1(w), f_2(w), \dots, f_{i-1}(w), f_{i+1}(w), \dots, f_n(w) \rangle$

▶ For each $i \in N$, $p_i : W \rightarrow [0, 1]$ is a probability measure.

For any formula φ , $[\varphi]$ is the set of states where φ is true

E.g.,
$$[s_i] = \{w \mid f_i(w) = s_i, [s_{-i}] = \{w \mid f_{-i}(w) = s_{-i}\}$$

For any φ , $w \in W$

$$p_{i,w}([\varphi]) = p_i([\varphi] \mid R_i(w)) = rac{p_i([\varphi] \cap R_i(w))}{p_i(R_i(w))}$$

Knowledge of Choice: For all $i \in N$, and all $w, v \in W$, if $w R_i v$, then $f_i(w) = f_i(v)$

Bayes Rational

Player *i* is **Bayes rational at** *w* if, for all $s_i \in S_i$,

$$\sum_{s_{-i} \in S_{-i}} p_{i,w}([s_{-i}]) \times u_i(f_i(w), s_{-i}) \ge \sum_{s_{-i} \in S_{-i}} p_{i,w}([s_{-i}]) \times u_i(s_i, s_{-i})$$

 $\varphi \Box \rightarrow \psi$ means "If it were the case that φ then it would be the case that ψ

 $s_i \square \rightarrow s_{-i}$ means "If it were the case that player *i* chose strategy s_i then it would be the case that her opponents chose strategy profile s_{-i}

 $\varphi \square \rightarrow \psi$ means "If it were the case that φ then it would be the case that ψ

 $s_i \square \rightarrow s_{-i}$ means "If it were the case that player *i* chose strategy s_i then it would be the case that her opponents chose strategy profile s_{-i}

Player *i* is **causally rational at** *w* if, for all $s_i \in S_i$,

$$\sum_{s_{-i}\in S_{-i}} p_{i,w}([f_i(w) \square \rightarrow s_{-i}]) \times u_i(f_i(w), s_{-i}) \ge$$
$$\sum_{s_{-i}\in S_{-i}} p_{i,w}([s_i \square \rightarrow s_{-i}]) \times u_i(s_i, s_{-i})$$

"If it were the case that player *i* choose strategy s_i , then it would be the case that her opponents choose strategy profile s_{-i} "

Lewis-Stalnaker Semantics

 $x \leq_w y$, x is "closer" to w than y.

For each $w \in W$, \leq_w is a relation on W such that for all $w \in W$,

- \leq_w is complete;
- \leq_w is transitive;
- \leq_w is antisymmetric;
- \leq_w is centered (for all $v, w \leq_w v$);

 $\min(X, w)$ is the minimal element in X at w according to \leq_w .

$$[\varphi \Box \rightarrow \psi] = \{w \mid \min([\varphi], w) \in [\psi]\}$$

$$p_{i,w}([arphi \square o \psi]) = rac{p_i([arphi \square o \psi] \cap R_i(w))}{p_i(R_i(w))}$$

Sufficiency For each player *i*, for every $s_i \in S_i$, there exists a state *w* such that $f_i(w) = s_i$.
Bob *L R U* 1,1 0,0 *D* 0,0 2,2







Eric Pacuit

$\begin{array}{c|c} & L & R \\ & U & 1,1 & 0,0 \\ D & 0,0 & 2,2 \end{array}$

Bob

Ann is Bayes rational at w_1 $p_{Ann,w_1}([L]) = 4/5, p_{Ann,w_1}([R]) = 1/5$ $4/5 \times 1 + 1/5 \times 0 > 4/5 \times 0 + 1/5 \times 2$





Eric Pacuit

Ann

Bob *L R U* 1,1 0,0 *D* 0,0 2,2

Ann is not causal rational at w_1 $[U \square L] = \{w_1\}, [U \square R] = \{w_2\}$ $4/5 \times 1 + 1/5 \times 0 = 0.8$





\leq_{w_1}	$\leq w_2$
W ₃	W 4
W2	W ₁
W4	W ₃
<i>w</i> ₁	W 2

Ann

Bob *L R U* 1,1 0,0 *D* 0,0 2,2

Ann is not causal rational at w_1 $[D \square L] = \{w_2\}, [D \square R] = \{w_1\}$ $2/5 \times 0 + 4/5 \times 2 = 1.6$



\leq_{w_1}	\leq_{W_2}
W ₃	W 4
W 2	W 1
W4	W ₃
<i>w</i> ₁	W2

Causal Independence

C1 For all *w* and *x*, for all $i \in N$, if $x \leq_w y$ for all *y* such that $f_i(y) = f_i(x)$, then $f_{-i}(x) = f_{-i}(w)$

Theorem. In any model G satisfying C1, player i is Bayes rational at w if and only if player i is causally rational at w.

O. Board. *The Equivalence of Bayes and Causal Rationality in Games*. Theory and Decision, 61, pgs. 1 - 19, 2006.

Theorem. In any model *G* satisfying C1, player *i* is Bayes rational at *w* if and only if player *i* is causally rational at *w*.

Proof. C1 implies that for all $s_i, s_{-i}, [s_i \square \rightarrow s_{-i}] = [s_{-i}]$

O. Board. *The Equivalence of Bayes and Causal Rationality in Games*. Theory and Decision, 61, pgs. 1 - 19, 2006.

- What lies behind the apparently paradoxical claim that a theory of rationality rests on a degree of fallibility on the part of the decision maker?
- How is ratifiability related to other theories of decision and to standard game-theoretic solution concepts?

H. S. Shin. A Reconstruction of Jeffrey's Notion of Ratifiability in Terms of Counterfactual Beliefs. Theory and Decision, 31, pgs. 21 - 47, 1991.

 $G = \langle S^1, S^2, u^1, u^2 \rangle$ is a two-player normal form game

- For $i = 1, 2, S^i$ is finite with K^i elements
- For $i = 1, 2, u^i : S^1 \times S^2 \rightarrow \mathbb{R}$
- For i = 1, 2, -i denotes player *i*'s opponent.

Let At be the following set of atomic propositions: for all $i = 1, 2, k = 1..., K^i$

- D_k^i means "player *i* decides to play s_k^i
- P_k^i means "player *i* performs s_k^i

A valuation is a function $V : At \rightarrow \{0, 1\}$. We say V is a **state** provided for all *i*, *k* and $j \neq k$,

$$V(D_k^i) = 1$$
 iff for all $j \neq k$, $V(D_j^i) = 0$
 $V(P_k^i) = 1$ iff for all $j \neq k$, $V(P_j^i) = 0$

Let Ω be the set of all states.

$$\delta_k^i = \{ V \mid V \in \Omega \text{ and } V(D_k^i) = 1 \}$$

$$\pi_k^i = \{ V \mid V \in \Omega \text{ and } V(P_k^i) = 1 \}$$

"By construction, these events do not coincide, and we leave open as a logical possibility the divergence between decisions and performances."

We have two partitions for each i = 1, 2:

•
$$\Delta^i = \{\delta^i_k \mid k = 1, \dots, K^i\}$$

$$\Pi^i = \{ \pi^i_k \mid k = 1, \dots, K^i \}$$

Let Δ be the meet of Δ^1 and Δ^2 and Π the meet of Π^1 and Π^2 .

Let *p* be a probability distribution over Ω .

 $U^{i}(k \mid j)$ is the expected utility of player *i* when she decides to play s_{k}^{i} but plays s_{i}^{i} instead.

$$U^{i}(k \mid j) := \sum_{m=1}^{K^{-i}} u^{i}(s_{k}^{i}, s_{m}^{-i})p(\pi_{m}^{-i} \mid \delta_{j}^{i} \cap \pi_{k}^{i})$$

 $(U^{i}(k \mid j) \text{ is only defined when } \delta^{i}_{i} \cap \pi^{i}_{k} \text{ is non-null under } p^{i})$

Let $\epsilon > 0$ be given. Assume $\epsilon < \min_i \{1/K^i\}$.

A1
$$p(\pi_j^i | \delta_k^i) = \epsilon$$
 for all $j \neq k$, whenever defined
A2 $p(\pi^i | \delta^i \cap \delta^{-i}) = p(\pi^i | \delta^i)$ for all π^i , δ^i , δ^{-i} whenever defined.
A3 $U^i(j | j) \ge U^j(k | j)$ for all j, k whenever defined.

p is ϵ -ratifiable for *i* if *p* satisfies A1, A2 and A3.

The precise *magnitude* of the trembles should play no part in the analysis. Rather, what matters is that such trembles exist, and that the be "small".

A4 $p(\pi_k^{-i}) = p(\delta_k^{-i})$ for all k

p is **modest for** *i* if it satisfies A4.

p is **modestly ratifiable for** *i* if there are sequences $p_1, p_2, ..., p_t, ...$ and $\epsilon_1, \epsilon_2, ..., \epsilon_t, ...$ such that for all *t*, p_t is modest for *i* and ϵ_t ratifiable for *i* and $p_t \rightarrow p$ as $\epsilon_t \rightarrow 0$.

Let *P* be a tremble-free distribution on Ω . I.e., $P(\delta_k^i) = P(\pi_k^i)$ for all *i*, *k*.

Define $\varphi(l \mid j) := P(\delta_l^{-i} \mid \delta_j^i)$ the probability that player *i*'s opponent (decides to) plays s_l^{-i} given that player *i* (decides to) play s_l^i .

 \overline{P} (a probability measure on $S^1 \times S^2$) is a **correlated equilibrium** provided for all *i*, *j*, *k* whenever $\varphi^i(I | j)$ is defined

$$\sum_{l=1}^{K^{-i}} \varphi^{i}(l \mid j) [u^{i}(s_{j}^{i}, s_{l}^{-i}) - u^{i}(s_{k}^{i}, s_{l}^{-i})] \geq 0$$



- Three Nash equilibria:
 - (*U*, *R*): the payoff is (2, 7)
 - (*D*, *L*): the payoff is (7, 2)
 - $([\frac{2}{3}(U), \frac{1}{3}D], [\frac{2}{3}(L), \frac{1}{3}(R)])$: the payoff is $(4\frac{2}{3}, 4\frac{2}{3})$



- Three Nash equilibria:
 - (*U*, *R*): the payoff is (2, 7)
 - (*D*, *L*): the payoff is (7, 2)
 - $([\frac{2}{3}(U), \frac{1}{3}D], [\frac{2}{3}(L), \frac{1}{3}(R)])$: the payoff is $(4\frac{2}{3}, 4\frac{2}{3})$
- After conducting the lottery, an outside observer provides Ann with a recommendation to play the first component of the profile that was chosen, and Bob the second component.



- Three Nash equilibria:
 - (*U*, *R*): the payoff is (2, 7)
 - (*D*, *L*): the payoff is (7, 2)
 - $([\frac{2}{3}(U), \frac{1}{3}D], [\frac{2}{3}(L), \frac{1}{3}(R)])$: the payoff is $(4\frac{2}{3}, 4\frac{2}{3})$
- After conducting the lottery, an outside observer provides Ann with a recommendation to play the first component of the profile that was chosen, and Bob the second component.
- ► The expected payoff is ¹/₃(6,6) + ¹/₃(2,7) + ¹/₃(7,2) = (5,5) (which is outside the convex hull of the Nash equilibria)

Theorem (Shin). p is modestly ratifiable if and only if p is a correlated equilibrium.

The distinction between *decisions* and *performances* is exactly analogous to the distinction between the *recommendations* issued by the arbitrator and the *actions* taken by the players.

Note that it is crucial that the players share the same probability measure p over the set of states.

When p is modestly ratifiable,

- No player will place positive probability on a strictly dominated strategy (So that in Newcomb's problem, both boxes are taken, and in the prisoners' dilemma, the players confess)
- In a two-person zero-sum game, the payoffs achievable by modest ratifiability cannot exceed the "value" of the game
- In non-zero sum games, the payoffs achievable by modest ratifiability can exceed the payoffs achieved as a Nash equilibrium.

Capture the notion of similarity by means of a *metric* on the space of possible worlds.

 $p \square \rightarrow q$ is true at *w* if and only if, there is a closed sphere *C* around *p* in the metric *m* such that $C \cap [p] \neq \emptyset$ and $C \cap [p] \subseteq [q]$

A possible world for player *i*, is a pair $\langle x, y \rangle$ where $x \in S^i$ and $y \in \Delta(S^{-i})$

Each state in an Aumann model can be associated with a *possible world* (for a fixed player)

Given $\langle W, \{R_i\}_{i \in N}, \{f_i\}_{i \in N}, \{p_i\}\rangle$, define $\beta^1 : W \to S^1 \times S$ (S is the unit simplex over the strategies of player *i*'s opponent)

$$\beta^{1}(w) = \langle f_{1}(w), \langle p_{i,w}([s^{-i}]) \rangle_{s^{-i} \in S^{-i}} \rangle$$







Eric Pacuit

"Library Stack Metric"

Let λ be a measure on player *i*'s possible worlds. For $\langle x, y \rangle, \langle x', y' \rangle \in \text{two}$ possible worlds for player *i*:

$$\lambda(\langle x, y \rangle, \langle x', y' \rangle) = \begin{cases} \sqrt{\sum_{j=1}^{2} |y_j - y'_j|^2} & \text{if } x = x' \\ \sqrt{\sum_{j=1}^{2} |y_j - y'_j|^2} + 1 & \text{if } x \neq x' \end{cases}$$













At β(w₁) = β(w₂), "If player 1 were to play D, his payoff would be higher" is false.





- At β(w₁) = β(w₂), "If player 1 were to play D, his payoff would be higher" is false.
- At β(w₃) = β(w₄), "If player 1 were to play U, his payoff would be higher" is false.


m-Rational

A player should never find himself at a possible world at which according to his metric m, he would be strictly better off if he were to deviate.

Let π_r^i mean "player *i*'s payoff does not exceed *r*" (for all $r \in \mathbb{R}$)

Suppose that $\beta^i(w) = \langle x, y \rangle$ and *m* is a metric in player *i*'s possible worlds. We say that *i* is *m*-rational at *w* if

For all $j \in \{1, ..., K^i\}$, there is a $r \leq U^i(\langle x, y \rangle)$ such that $\sigma_j^i \square \pi_r^i$ is true at $\langle x, y \rangle$.

m-Rationality describes a type of consistency—the consistency of a player's probability distribution with his metric.

Theorem. Player *i* is λ -rational at *w* if and only if *i* is Aumann rational at *w*.

Corollary. Suppose that $p^i = p$ for all *i*. Then ally players are λ -rational if and only if *p* is a correlated equilibrium.

Corollary. Suppose that $p^i = p$ for all *i* and *p* is independent. Then all players are λ -rational if and only if the mixed strategies given by *p* is a Nash equilibrium.