Clear Thinking in an Uncertain World: Human Reasoning and its Foundations Lecture 12

Eric Pacuit

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

November 18, 2013

Two Puzzles about Rationality and Coordination

- 1. The Prisoner's Dilemma
- 2. Newcomb's Paradox

L R1 0



L R U 11 00 D 00 11







What should Ann (Bob) do?



What should Ann (Bob) do?

Just Enough Game Theory

"Game theory is a bag of analytical tools designed to help us understand the phenomena that we observe when decision-makers interact."

Osborne and Rubinstein. Introduction to Game Theory. MIT Press .

Just Enough Game Theory

"Game theory is a bag of analytical tools designed to help us understand the phenomena that we observe when decision-makers interact."

Osborne and Rubinstein. Introduction to Game Theory. MIT Press .

A game is a description of strategic interaction that includes

- actions the players can take
- description of the players' interests (i.e., preferences),
- description of the "structure" of the decision problem

Just Enough Game Theory

"Game theory is a bag of analytical tools designed to help us understand the phenomena that we observe when decision-makers interact."

Osborne and Rubinstein. Introduction to Game Theory. MIT Press .

A game is a description of strategic interaction that includes

- actions the players can take
- description of the players' interests (i.e., preferences),
- description of the "structure" of the decision problem

It does not specify the actions that the players do take.

A **solution concept** is a systematic description of the outcomes that may emerge in a family of games.

This is the starting point for most of game theory and includes many variants: Nash equilibrium, backwards inductions, or iterated dominance of various kinds.

These are usually thought of as the embodiment of "rational behavior" in some way and used to analyze game situations.

A **solution concept** is a systematic description of the outcomes that may emerge in a family of games.

This is the starting point for most of game theory and includes many variants: Nash equilibrium, backwards inductions, or iterated dominance of various kinds.

These are usually thought of as the embodiment of "rational behavior" in some way and used to analyze game situations.

For this course, solution concepts are more of an endpoint.





Clear Thinking in an Uncertain World







What should Ann do?



What should Ann do? Bob best choice in Ann's worst choice



What should Ann do? maximize over each row and choose the maximum value



What should Bob *do*? *minimize over each column and choose the maximum value*



Von Neumann Minmax Theorem. In any finite, two-player, zero-sum game, there is always at least one minmax solution.



What is a rational choice for Ann (Bob)?



What is a rational choice for Ann (Bob)? Flip a coin!



What is a rational choice for Ann (Bob)?



What is a rational choice for Ann (Bob)? Play a different game!

Two people commit a crime.

Two people commit a crime. The are arrested by the police, who are quite sure they are guilty but cannot prove it without at least one of them confessing.

Two people commit a crime. The are arrested by the police, who are quite sure they are guilty but cannot prove it without at least one of them confessing. The police offer the following deal. Each one of them can confess and get credit for it.

Two people commit a crime. The are arrested by the police, who are quite sure they are guilty but cannot prove it without at least one of them confessing. The police offer the following deal. Each one of them can confess and get credit for it. If only one confesses, he becomes a state witness and not only is he not punished, he gets a reward.

Two people commit a crime. The are arrested by the police, who are quite sure they are guilty but cannot prove it without at least one of them confessing. The police offer the following deal. Each one of them can confess and get credit for it. If only one confesses, he becomes a state witness and not only is he not punished, he gets a reward. If both confess, they will be punished but will get reduced sentences for helping the police.

Two people commit a crime. The are arrested by the police, who are quite sure they are guilty but cannot prove it without at least one of them confessing. The police offer the following deal. Each one of them can confess and get credit for it. If only one confesses, he becomes a state witness and not only is he not punished, he gets a reward. If both confess, they will be punished but will get reduced sentences for helping the police. If neither confesses, the police honestly admit that there is no way to convict them, and they are set free.

Two options: Confess (C), Don't Confess (D)
Two options: Confess (C), Don't Confess (D)

Possible outcomes:

Two options: Confess (C), Don't Confess (D)

Possible outcomes: We both confess (C, C),

Two options: Confess (C), Don't Confess (D)

Possible outcomes: We both confess (C, C), I confess but my partner doesn't (C, D),

Two options: Confess (C), Don't Confess (D)

Possible outcomes: We both confess (C, C), I confess but my partner doesn't (C, D), My partner confesses but I don't (D, C),

Two options: Confess (C), Don't Confess (D)

Possible outcomes: We both confess (C, C), I confess but my partner doesn't (C, D), My partner confesses but I don't (D, C), neither of us confess (D, D).





Ann's preferences

Clear Thinking in an Uncertain World



Bob's preferences

Clear Thinking in an Uncertain World



Dominance Reasoning



Dominance Reasoning



Dominance Reasoning







What should Ann (Bob) do? Dominance reasoning



What should Ann (Bob) do? Dominance reasoning



What should Ann (Bob) do? Dominance reasoning is not Pareto!



What should Ann (Bob) do? Think as a group!



What should Ann (Bob) do? Play against your mirror image!



What should Ann (Bob) do? Play against your mirror image!



What should Ann (Bob) do? *Change the game* (eg., Symbolic Utilities)





Assurance Game



"Yet the symbolic value of an act is not determined solely by *that* act.

"Yet the symbolic value of an act is not determined solely by *that* act. The act's meaning can depend upon what other acts are available with what payoffs and what acts also are available to the other party or parties.

"Yet the symbolic value of an act is not determined solely by *that* act. The act's meaning can depend upon what other acts are available with what payoffs and what acts also are available to the other party or parties. What the act symbolizes is something it symbolizes when done in *that* particular situation, in preference to *those* particular alternatives.

"Yet the symbolic value of an act is not determined solely by *that* act. The act's meaning can depend upon what other acts are available with what payoffs and what acts also are available to the other party or parties. What the act symbolizes is something it symbolizes when done in *that* particular situation, in preference to *those* particular alternatives. If an act symbolizes "being a cooperative person," it will have that meaning not simply because it has the two possible payoffs it does

"Yet the symbolic value of an act is not determined solely by *that* act. The act's meaning can depend upon what other acts are available with what payoffs and what acts also are available to the other party or parties. What the act symbolizes is something it symbolizes when done in *that* particular situation, in preference to those particular alternatives. If an act symbolizes "being a cooperative person," it will have that meaning not simply because it has the two possible payoffs it does but also because it occupies a particular position within the two-person matrix — that is, being a dominated action that (when joined with the other person's dominated action) yield a higher payoff to each than does the combination of dominated actions. " (pg. 55)

R. Nozick. The Nature of Rationality. Princeton University Press, 1993.



Prisoner's Dilemma
















What should/will Ann (Bob) do?

"Game theorists think it just plain wrong to claim that the Prisoners' Dilemma embodies the essence of the problem of human cooperation. "Game theorists think it just plain wrong to claim that the Prisoners' Dilemma embodies the essence of the problem of human cooperation. On the contrary, it represents a situation in which the dice are as loaded against the emergence of cooperation as they could possibly be. If the great game of life played by the human species were the Prisoner's Dilemma, we wouldn't have evolved as social animals! "Game theorists think it just plain wrong to claim that the Prisoners' Dilemma embodies the essence of the problem of human cooperation. On the contrary, it represents a situation in which the dice are as loaded against the emergence of cooperation as they could possibly be. If the great game of life played by the human species were the Prisoner's Dilemma, we wouldn't have evolved as social animals! No paradox of rationality exists. Rational players don't cooperate in the Prisoners' Dilemma, because the conditions necessary for rational cooperation are absent in this game." (pg. 63)

K. Binmore. Natural Justice. Oxford University Press, 2005.

Two boxes in front of you, A and B.

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

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).



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.

	B = 1M	B = 0
1 Box	1M	0
2 Boxes	1M + 1000	1000



	B = 1M	B = 0		
1 Box	1M	0		
2 Boxes	1M + 1000	1000	Ī	

	B = 1M	B = 0
1 Box	h	1-h
2 Boxes	1-h	h



J. Collins. *Newcomb's Problem*. International Encyclopedia of Social and Behavorial Sciences, 1999.

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*.

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*. (A nasty nephew wants inheritance from his rich Aunt.

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*. (A nasty nephew wants inheritance from his rich Aunt. The nephew wants the inheritance, but other things being equal, does not want to apologize.

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*. (A nasty nephew wants inheritance from his rich Aunt. The nephew wants the inheritance, but other things being equal, does not want to apologize. Does dominance give the nephew a reason to not apologize?

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*. (A nasty nephew wants inheritance from his rich Aunt. The nephew wants the inheritance, but other things being equal, does not want to apologize. Does dominance give the nephew a reason to not apologize? Whether or not the nephew is cut from the will may depend on whether or not he apologizes.)

There is a conflict between maximizing your expected value (1-box choice) and dominance reasoning (2-box choice).

Dominance reasoning is appropriate only when probability of outcome is *independent of choice*. (A nasty nephew wants inheritance from his rich Aunt. The nephew wants the inheritance, but other things being equal, does not want to apologize. Does dominance give the nephew a reason to not apologize? Whether or not the nephew is cut from the will may depend on whether or not he apologizes.)

What the Predictor did yesterday is *probabilistically dependent* on the choice today, but *causally independent* of today's choice.

 $V(A) = \sum_{w} V(w) \cdot P_A(w)$

(the expected value of act A is a probability weighted average of the values of the ways w in which A might turn out to be true)

 $V(A) = \sum_{w} V(w) \cdot P_A(w)$ (the expected value of act A is a probability weighted average of the values of the ways w in which A might turn out to be true)

Orthodox Bayesian Decision Theory: $P_A(w) := P(w \mid A)$ (Probability of w given A is chosen)

Causal Decision theory: $P_A(w) = P(A \Box \rightarrow w)$ (Probability of *if A* were chosen then w would be true)

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

```
V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1)
```

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

```
V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1) = 1000000 \cdot 0.99 + 0 \cdot 0.01
```

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

```
V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1) = 1000000 \cdot 0.99 + 0 \cdot 0.01 = 990,000
```

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

$$V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1) = 1000000 \cdot 0.99 + 0 \cdot 0.01 = 990,000$$

$$V(B_2) = V(L)P(L \mid B_2) + V(K)P(K \mid B_2)$$

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

```
V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1) = 1000000 \cdot 0.99 + 0 \cdot 0.01 = 990,000
```

```
V(B_2) = V(L)P(L \mid B_2) + V(K)P(K \mid B_2) = 1001000 \cdot 0.01 + 1000 \cdot 0.99
```

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

```
V(B_1) = V(M)P(M \mid B_1) + V(N)P(N \mid B_1) = 1000000 \cdot 0.99 + 0 \cdot 0.01 = 990,000
```

$$V(B_2) = V(L)P(L \mid B_2) + V(K)P(K \mid B_2) =$$

1001000 · 0.01 + 1000 · 0.99 = 11,000

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- *M*: receive \$1,000,000
- L: receive \$1,001,000

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N)$

B1: one-box (open box B)
B2: two-box choice (open both A and B)
N: receive nothing
K: receive \$1,000
M: receive \$1,000,000
L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N) = 1000000 \cdot \mu + 0 \cdot 1 - \mu$

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- K: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N) = 1000000 \cdot \mu + 0 \cdot 1 - \mu = 1000000 \mu$

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- *K*: receive \$1,000
- M: receive \$1,000,000
- L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N) = 1000000 \cdot \mu + 0 \cdot 1 - \mu = 1000000 \mu$

$$V(B_2) = V(L)P(B_2 \Box \rightarrow L) + V(K)P(B_2 \Box \rightarrow K)$$

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- *K*: receive \$1,000
- *M*: receive \$1,000,000
- L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N) = 1000000 \cdot \mu + 0 \cdot 1 - \mu = 1000000\mu$

$$V(B_2) = V(L)P(B_2 \Box \rightarrow L) + V(K)P(B_2 \Box \rightarrow K) =$$

1001000 · μ + 1000 · 1 - μ

- B_1 : one-box (open box B)
- B_2 : two-box choice (open both A and B)
- N: receive nothing
- *K*: receive \$1,000
- *M*: receive \$1,000,000
- L: receive \$1,001,000

 $V(B_1) = V(M)P(B_1 \Box \rightarrow M) + V(N)P(B_1 \Box \rightarrow N) = 1000000 \cdot \mu + 0 \cdot 1 - \mu = 1000000\mu$

$$V(B_2) = V(L)P(B_2 \Box \to L) + V(K)P(B_2 \Box \to K) = 1001000 \cdot \mu + 1000 \cdot 1 - \mu = 1000000\mu + 1000$$

D. Lewis. *Prisoner's Dilemma Is a Newcomb Problem*. Philosophy and Public Affairs, 8, pgs. 235-240, 1979.

S. Brams. *Newcomb's Problem and Prisoners' Dilemma*. The Journal of Conflict Resolution, 19:4, pgs. 596 - 612, 1975.

S. Hurley. *Newcomb's Problem, Prisoner's Dilemma and Collective Action*. Synthese 86, pgs. 173 - 196, 1991.