



DEPARTMENT OF ECONOMICS
Working Paper Series

Counterfactuals and the Prisoner's Dilemma

Giacomo Bonanno
U.C. Davis

May 17, 2013

Paper # 13-7

This is the first draft of a chapter in a planned book on the Prisoner's Dilemma, edited by Martin Peterson, to be published by Cambridge University Press. It discusses the nature of the conditionals involved in deliberation, taking the Prisoner's Dilemma game as point of departure.

Department of Economics
One Shields Avenue
Davis, CA 95616
(530)752-0741

http://www.econ.ucdavis.edu/working_search.cfm

Counterfactuals and the Prisoner's Dilemma

Giacomo Bonanno*

Department of Economics
University of California
Davis, CA 95616 – USA
e-mail: gfbonanno@ucdavis.edu

1. Introduction

In 2011 Harold Camping, president of Family Radio (a California-based Christian radio station), predicted that Rapture (the taking up into heaven of God's elect people) would take place on May 21, 2011. In light of this prediction some of his followers gave up their jobs, sold their homes and spent large sums promoting Camping's claims. Did these people act rationally? Consider also the following hypothetical scenarios. Early in 2012, on the basis of a popular reading of the Mayan calendar, Ann believes that the world is going end on December 21, 2012. She drops out of college, withdraws all the money she has in her bank account and decides to spend it all on travelling and enjoying herself. Is her decision rational? Bob smokes two packets of cigarettes a day; when asked if he would still smoke if he knew that he was going to get lung cancer from smoking, he answers 'No'; when asked if he is worried about getting lung cancer, he says that he is not and explains that his grandfather was a heavy smoker all his life and died – cancer free – at the age of 98. Bob believes that, like his grandfather, he is immune from lung cancer. Is Bob's decision to continue smoking rational?

I will argue below that the above questions are closely related to the issue, hotly debated in the literature, whether it can be rational for the players to choose "Cooperation" in the Prisoner's Dilemma game, shown in Figure 1. It is a two-player, simultaneous game where each player has two strategies: "Cooperation" (denoted by C for Player 1 and by c for Player 2) and "Defection" (denoted by D for Player 1 and by d for Player 2). Part *a* of Figure 1 shows the outcomes associated with each strategy-pair. Player 1's ranking of the outcomes is $z_3 \succ_1 z_1 \succ_1 z_4 \succ_1 z_2$ (the interpretation of $x \succ_1 y$ is that Player 1 strictly prefers outcome x to

outcome y) and Player 2's ranking is $z_2 \succ_2 z_1 \succ_2 z_4 \succ_2 z_3$. In Part *b* of Figure 1 ordinal utility functions are used to represent the players' ranking. An ordinal utility function $U_1 : \{z_1, z_2, z_3, z_4\} \rightarrow \mathbb{R}$ (where \mathbb{R} denotes the set of real numbers) is said to represent Player 1's ranking if it satisfies the property that, for any two outcomes x and y , $U_1(x) > U_1(y)$ if and only if $x \succ_1 y$. Similarly for player 2. In each cell of the table, the first number is the utility of Player 1 and the second number the utility of Player 2.¹

		Player 2			
		<i>c</i>	<i>d</i>		
Player 1	<i>C</i>	z_1	z_2	Player 1	<i>C</i>
	<i>D</i>	z_3	z_4		<i>D</i>

a

Player 1's ranking of the outcomes is
 $z_3 \succ_1 z_1 \succ_1 z_4 \succ_1 z_2$
 Player 2's ranking of the outcomes is
 $z_2 \succ_2 z_1 \succ_2 z_4 \succ_2 z_3$

		Player 2			
		<i>c</i>	<i>d</i>		
Player 1	<i>C</i>	2, 2	0, 3	Player 1	<i>C</i>
	<i>D</i>	3, 0	1, 1		<i>D</i>

b

Ordinal utility functions used to represent the rankings

Figure 1
The Prisoner's Dilemma game

What constitutes a rational choice for a player? We take the following to be the basic definition of rationality:

* This is the first draft of a chapter in a planned book on the Prisoner's Dilemma, edited by Martin Peterson, to be published by Cambridge University Press.

¹ We take rankings of (that is, preferences over) the outcomes as primitives (and utility functions merely as tools for representing those rankings). Thus we are not following the *revealed preference* approach, where observed choices are the primitives and preferences (or utility) are a derived notion:

“In revealed-preference theory, it isn't true [...] that Pandora chooses *b* rather than *a* because she prefers *b* to *a*. On the contrary, it is because Pandora chooses *b* rather than *a* that we say that Pandora prefers *b* to *a*, and assign *b* a larger utility.” (Binmore, 2011, p. 19.)

Thus in the Prisoner's Dilemma game of Part *b* of Figure 1,

“Writing a larger payoff for Player 1 in the bottom-left cell of the payoff table than in the top-left cell is just another way of registering that Player 1 would choose *D* if she knew that Player 2 were going to choose *c*. [W]e must remember that Player 1 doesn't choose *D* because she then gets a larger payoff. Player 1 assigns a larger payoff to [the outcome associated with] (*D*,*c*) than to [the outcome associated with] (*C*,*c*) because she would choose the former if given the choice.” (Binmore, 2011, pp. 27-28, with minor modifications to adapt the quotation to the notation used in Figure 1.)

For a criticism of (various interpretations of) the notion of revealed preference see Chapter 3 of Hausman, 2012; see also Rubinstein and Salant, 2008.

*A choice is rational if it is optimal given the
decision-maker's preferences and beliefs.* (BDR)

More precisely, we say that it is rational for the decision-maker to choose action *a* if there is no other feasible action *b* which – *according to her beliefs* – would yield an outcome that she prefers to the outcome that – again, according to her beliefs – would be a consequence of taking action *a*. According to this definition, the followers of Harold Camping did act rationally when they decided to sell everything and devote themselves to promoting Camping's claim: they believed that the world was soon coming to an end and – presumably – they viewed their proselytizing as “qualifying them for Rapture”, undoubtedly an outcome they preferred to the alternative of enduring the wrath of Judgment Day. Similarly, Ann's decision to live it up in anticipation of the end of the world predicted by the Mayan calendar qualifies as rational, as does Bob's decision to carry on smoking on the belief that – like his grandfather – he will be immune from cancer. Thus anybody who argues that the above decisions are *not* rational must be appealing to a stronger definition of rationality than *BDR*: for example, one could question the rationality of holding those beliefs.

When the rationality of beliefs is called into question, an asymmetry is introduced between preferences and beliefs. Concerning preferences it is a generally accepted principle that *de gustibus non est disputandum* (in matters of taste, there can be no disputes). According to this principle, there is no such thing as an irrational preference. As Rubinstein notes,

“According to the assumption of rationality in economics, the decision maker is guided by his preferences. But the assumption does not impose a limitation on the reasonableness of preferences. The preferences can be even in direct contrast with what common sense might define as the decision maker's interests.” (Rubinstein, 2012, p. 49.)

For example, I cannot be judged to be irrational if I prefer an immediate benefit (e.g. from taking a drug) with known negative future consequences (e.g. from addiction) over an immediate sacrifice (e.g. by enduring pain) followed by better long-term health.²

In the matter of beliefs, on the other hand, it is generally thought that one *can* contend that some particular beliefs are “unreasonable” or “irrational”, by appealing to such arguments as

² For a criticism of the view that preferences are not subject to rational scrutiny see Chapter 10 of Hausman, 2012.

the lack of supporting evidence, the incorrect processing of relevant information, the denial of laws of Nature, etc.

Consider now the following statement by Player 1 in the Prisoner's Dilemma ('COR' stands for 'correlation'):

“I believe that if I play C then Player 2 will play c and that if I play D then Player 2 will play d . Hence, if I play C then the outcome will be z_1 and if I play D then the outcome will be z_4 . Since I prefer z_1 to z_4 , I have decided to play C .” (COR₁)

Given the reported beliefs, Player 1's decision to play C is rational according to definition *BDR*. Thus in order to question the rationality of Player 1's decision, one has to argue that the beliefs expressed in COR₁ violate some principle of rationality. In the literature, there are those who claim that Player 1's reported beliefs are irrational and those who claim that those beliefs can be rationally justified, for example by appealing to the symmetry of the game (see, for example, Brams, 1975, and Davis, 1977, 1985) or to special circumstances, such as the players being identical in some sense (e.g. they are identical twins): this has become known as the “Identity Assumption” (see Bicchieri and Green, 1999, and Gilboa, 1999).

In order to elucidate what is involved in Player 1's belief “if I play C then Player 2 will play c , and if I play D then Player 2 will d ” we broaden the discussion to the more general issue of the role of beliefs and conditionals in game-theoretic reasoning.

2. Models of games: beliefs and choices

It is a widely held opinion that the notion of rationality involves the use of counterfactual reasoning. For example, Aumann writes:

“[O]ne really cannot discuss rationality, or indeed decision making, without substantive conditionals and counterfactuals. Making a decision means choosing among alternatives. Thus one must consider hypothetical situations – what would happen if one did something different from what one actually does. [I]n interactive decision making – games – you must consider what other people would do if you did something different from what you actually do.” (Aumann, 1995, p. 15.)

How is counterfactual reasoning incorporated in the analysis of games? The definition of strategic-form game provides only a partial description of an interactive situation. A *game in strategic form with ordinal payoffs* is defined as a quintuple $G = \langle N, \{S_i\}_{i \in N}, O, z, \{\succeq_i\}_{i \in N} \rangle$, where $N = \{1, \dots, n\}$ is a set of *players*, S_i is the set of *strategies* of (or possible choices for) player $i \in N$, O is a set of possible *outcomes*, $z : S \rightarrow O$ is a function that associates an outcome with every strategy profile $s = (s_1, \dots, s_n) \in S = S_1 \times \dots \times S_n$ and \succeq_i is a complete and transitive binary relation on O representing player i 's ranking of the outcomes (the interpretation of $o \succeq_i o'$ is that player i considers outcome o to be at least as good as outcome o'). Games are typically represented in *reduced form* by replacing the triple $\langle O, z, \{\succeq_i\}_{i \in N} \rangle$ with a set of *payoff functions* $\{\pi_i\}_{i \in N}$, where $\pi_i : S \rightarrow \mathbb{R}$ is any numerical function that satisfies the property that, for all $s, s' \in S$, $\pi_i(s) \geq \pi_i(s')$ if and only if $z(s) \succeq_i z(s')$, that is, if player i considers the outcome associated with s to be at least as good as the outcome associated with s' . In the following we will adopt this more succinct representation of strategic-form games.³ Thus the definition of strategic-form game only specifies what choices each player has available and how the player ranks the possible outcomes; it is silent on what the player believes. In order to complete the description one needs to introduce the notion of *model of a game*.

Definition 1. Given a strategic-form game G , a *model of G* is a triple $\langle \Omega, \{\sigma_i\}_{i \in N}, \{\mathcal{B}_i\}_{i \in N} \rangle$ where Ω is a set of *states* and, for every player $i \in N$, $\sigma_i : \Omega \rightarrow S_i$ is a function that associates with every state $\omega \in \Omega$ a strategy $\sigma_i(\omega) \in S_i$ of player i and $\mathcal{B}_i \subseteq \Omega \times \Omega$ is a binary relation representing the beliefs of player i . The interpretation of $\omega \mathcal{B}_i \omega'$ is that at state ω player i considers state ω' possible. Let $\mathcal{B}_i(\omega) = \{\omega' \in \Omega : \omega \mathcal{B}_i \omega'\}$; thus $\mathcal{B}_i(\omega)$ is the set of states that

³ It is important to note, however, that the payoff functions are taken to be purely ordinal and one could replace π_i with any other function obtained by composing π_i with an arbitrary strictly increasing function on the set of real numbers. In the literature it is customary to impose a stronger assumption on players' preferences, namely that each player has a complete and transitive preference relation on the set of probability distributions over the set of outcomes O , which satisfies the axioms of Expected Utility. For our purposes this stronger assumption is not needed.

player i considers possible at state ω .⁴

The functions $\{\sigma_i : \Omega \rightarrow S_i\}_{i \in N}$ give content to the players' beliefs. If $\sigma_i(\omega) = x \in S_i$ then the usual interpretation is that at state ω player i “chooses” strategy x . The exact meaning of ‘choosing’ is not elaborated further in the literature: does it mean that player i *has actually played* x or that she is *committed to playing* x or that x is the *output of her deliberation process*? We will adopt the latter interpretation: ‘player i chooses x ’ will be taken to mean ‘player i has made up her mind to play x ’.

Part *a* of Figure 2 shows a strategic-form game in reduced form and Part *b* a model of it. We represent the relation \mathcal{B}_i graphically as follows: $\omega \mathcal{B}_i \omega'$ (or, equivalently, $\omega' \in \mathcal{B}_i(\omega)$) if and only if there is an arrow from ω to ω' .

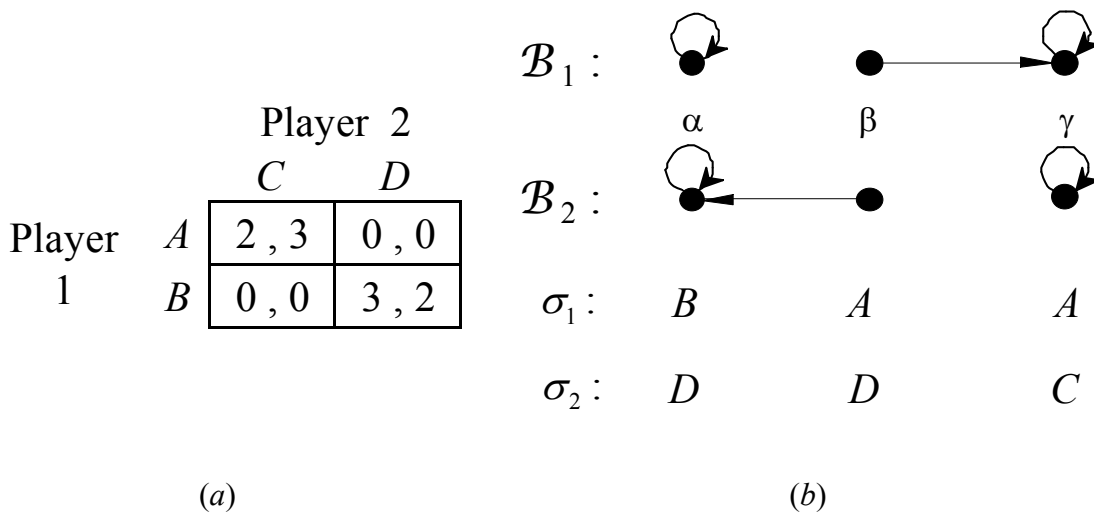


Figure 2

(a) A strategic-form game in reduced form. (b) A model of the game.

State β in the model of Figure 2 represents the following situation: Player 1 has made up her mind to play A and Player 2 has made up his mind to play D ; Player 1 erroneously believes that

⁴ Thus the relation \mathcal{B}_i can also be viewed as a function $\mathcal{B}_i : \Omega \rightarrow 2^\Omega$; such functions are called *possibility correspondences* in the literature. For further details on the so called “epistemic foundation program” in game theory, the reader is referred to Battigalli and Bonanno, 1999.

Player 2 has made up his mind to play C ($\mathcal{B}_1(\beta) = \{\gamma\}$ and $\sigma_2(\gamma) = C$) and Player 2 erroneously believes that Player 1 has made up her mind to play B ($\mathcal{B}_2(\beta) = \{\alpha\}$ and $\sigma_1(\alpha) = B$).

Remark 1. The model of Figure 2 reflects a standard assumption in the literature, namely that *a player is never uncertain about her own choice*: any uncertainty has to do with the other players' choices. This requirement is expressed formally as follows: for every $\omega' \in \mathcal{B}_i(\omega)$, $\sigma_i(\omega') = \sigma_i(\omega)$. We shall revisit this point in Section 5.

Returning to the model of Part *b* of Figure 2, a natural question to ask is whether the players are rational at state β . Consider Player 1: according to her beliefs, the outcome is going to be the one associated with the strategy pair (A, C) , with a corresponding payoff of 2 for her. In order to determine whether the decision to play A is rational, Player 1 needs to answer the question “what would happen if, instead of playing A , I were to play B ?”. The model is silent about such counterfactual scenarios. Thus the definition of model introduced above appears to lack the resources to address the issue of rational choice. Before we discuss how to enrich the definition of model, we turn, in Section 3, to a brief digression on the notion of counterfactual.

It should be noted, however, that a large literature – that originates in Aumann, 1987 – develops the analysis of rationality in strategic-form games using the models described above, without enriching them with an explicit framework for counterfactuals. However, as Shin (1992, p. 412) notes “If counterfactuals are not explicitly invoked, it is because the assumptions are buried implicitly in the discussion.” We shall return to this point in Section 4.

3. Stalnaker-Lewis selection functions

There are different types of conditionals. A conditional of the form “If John received my message he will be here soon” is called an *indicative* conditional. Conditionals of the form “If I were to drop this vase, it would break” and “If we had not missed the connection, we would be at home now” are called *subjunctive* conditionals; the latter is also an example of a *counterfactual*, namely a conditional with a false antecedent (we did in fact miss the connection). It is controversial how best to classify conditionals and we will not address this issue here. We are interested in the use of conditionals in the analysis of games and thus the relevant conditionals are those that pertain to deliberation.

In the decision-theoretic and game-theoretic literature the conditionals involved in deliberation are usually called “counterfactuals”, as illustrated in the quotation from Aumann (1995) in Section 2 and in the following:

“[R]ational decision-making involves conditional propositions: when a person weighs a major decision, it is rational for him to ask, for each act he considers, what would happen if he performed that act. It is rational, then, for him to consider propositions of the form ‘If I were to do a , then c would happen’. Such a proposition we shall call a counterfactual.” (Gibbard and Harper, 1978, p. 153.)

With the exception of Shin (1992), Bicchieri and Green (1999), Zambrano (2004) and Board (2006) (whose contributions are discussed in Section 4), the issue of counterfactual reasoning in strategic-form games has not been dealt with explicitly in the literature.⁵

We denote by $\phi > \psi$ the conditional “if ϕ were the case then ψ would be the case”. In the Stalnaker-Lewis theory of conditionals (Stalnaker, 1968, Lewis, 1973) the formula $\phi > \psi$ has a truth value which is determined as follows: $\phi > \psi$ is true at a state ω if and only if ψ is true at all the ϕ -states ω' that are closest (that is, most similar) to ω (a state ω' is a ϕ -state if ϕ is true at ω'). While Stalnaker postulates that, for every state ω and formula ϕ there is a unique ϕ -state ω' that is closest to ω , Lewis allows for the possibility that there may be several such states.

The semantic representation of conditionals is done by means of a *selection function* $f : \Omega \times 2^\Omega \rightarrow 2^\Omega$ (where 2^Ω denotes the set of subsets of Ω) that associates with every state ω and subset $E \subseteq \Omega$ (representing a proposition) a subset $f(\omega, E) \subseteq E$ interpreted as the states in E that are closest to ω . Several restrictions are imposed on the selection function, but we will skip the details.⁶

⁵ Although the issue has been discussed extensively in the context of dynamic games. See Bonanno (2013a) for a general discussion and relevant references.

⁶ For example, the restriction that if $\omega \in E$ then $f(\omega, E) = \{\omega\}$.

Just as the notion of accessibility relation enables us to represent a player's beliefs without, in general, imposing any restrictions on the content of those beliefs, the notion of selection function enables us to incorporate subjunctive conditionals into a model without imposing any constraints on what ϕ -states ought to be considered most similar to a state where ϕ is not true. A comic strip on the web site <http://xkcd.com/1170/> shows the following dialogue between father and son:

Father: No, you can't go.

Son: But all my friends ...

Father: If all your friends jumped off a bridge, would you jump too?

Son: Oh, Jeez... Probably.

Father: What!? Why!?

Son: Because all my friends did. Think about it: which scenario is more likely? Every single friend I know – many of them levelheaded and afraid of heights – abruptly went crazy at exactly the same time ...or the bridge is on fire?

The issue of determining what state(s) ought to be deemed closest to a given state is not a straightforward one. Usually “closeness” is interpreted in terms of a *ceteris paribus* (other things being equal) condition. However, typically *some* background conditions *must* be changed in order to evaluate a counterfactual. Consider, for example, the situation represented by state β in the model of Figure 2. What would be – in an appropriately enriched model – the closest state to β , call it η , where Player 1 plays B rather than A ? It has been argued (we will return to this point later) that it ought to be postulated that η is a state where Player 1 has the same beliefs about Player 2's choice as in state β . But, if – given Player 1's beliefs at β – the choice of A is rational, then at η one of the background conditions that describe state β no longer holds, namely, that Player 1 is rational and knows that she is rational. Alternatively, if one wants to hold this condition constant, then one must postulate that at η Player 1 believes that Player 2 is playing D and thus one must change another background condition at β , namely her beliefs about Player 2. We will return to these issue in Section 4.

There is also another issue that needs to be addressed. The selection function f is usually interpreted as capturing the notion of “causality” or “objective counterfactuality”. For example, suppose that Ann is facing two faucets, one labeled ‘hot’ and the other ‘cold’, and she needs hot water. Suppose also that the faucets are mislabeled and Ann is unaware of this. Then it would be objectively or causally true that “if Ann turned on the faucet labeled ‘cold’ she would get hot

water”; however, she could not be judged to be irrational if she expressed the belief “if I turned on the faucet labeled ‘cold’ I would get cold water” (and acted on this belief by turning on the faucet labeled ‘hot’). Since what we are interested in is the issue of rational choice, objective counterfactuals do not seem to be the relevant objects: *what matters is not what would in fact be the case but what the agent believes would be the case*. We shall call such beliefs *subjective counterfactuals*. How should these subjective counterfactuals be modeled? There are two options, examined in the following sub-sections.

3.1 Subjective counterfactuals as dispositional belief revision

One construal of subjective counterfactuals is in terms of a *subjective selection function* $f_i : \Omega \times 2^\Omega \rightarrow 2^\Omega$ such that, for every $\omega \in \Omega$ and $E \subseteq \Omega$, $f_i(\omega, E) \subseteq E$. The function f_i is interpreted as expressing, at every state, player i ’s *initial beliefs together with her belief revision policy*. Fix a state $\omega \in \Omega$ and consider the function $f_{i,\omega} : 2^\Omega \rightarrow 2^\Omega$ given by $f_{i,\omega}(E) = f_i(\omega, E)$, for every $E \subseteq \Omega$. Then this function gives the initial beliefs of player i at state ω (represented by the set $f_{i,\omega}(\Omega)$) as well as the set of states that player i would consider possible, at state ω , under the supposition that event $E \subseteq \Omega$ is true (represented by the set $f_{i,\omega}(E)$), for every event E . Subjective selection functions – with the implied dispositional belief revision policy – have been used extensively in the literature on dynamic games,⁷ but (to the best of my knowledge) have not been used in the analysis of strategic-form games, with the exception of Shin (1992) and Zambrano (2004), whose contributions are discussed in Section 4.

In this context, an enriched model of a strategic-form game G is a quadruple $\langle \Omega, \{\sigma_i\}_{i \in N}, \{\mathcal{B}_i\}_{i \in N}, \{f_i\}_{i \in N} \rangle$, where $\langle \Omega, \{\sigma_i\}_{i \in N}, \{\mathcal{B}_i\}_{i \in N} \rangle$ is as defined in Definition 1 and, for every player i , $f_i : \Omega \times 2^\Omega \rightarrow 2^\Omega$ is a subjective selection function satisfying the property that, for

⁷ See, for example, Arló-Costa and Bicchieri (2007), Board (2004), Clausen (2004), Halpern (1999, 2001), Rabinowicz (2000), Stalnaker (1998). For a critical discussion of this approach see Bonanno (2103a).

every state ω , $f_i(\omega, \Omega) = \mathcal{B}_i(\omega)$.⁸ Such enriched models would be able to capture the following reasoning of Player 1 in the Prisoner's Dilemma (essentially a restatement of COR_1):

“I have chosen to play C and I believe that Player 2 has chosen to play c and thus I believe that my payoff will be 2; furthermore, I am happy with my choice of C because – under the supposition that I play D – I believe that my payoff would be 1.” (COR_2)

These beliefs are illustrated by state α in the following enriched model of the Prisoner's Dilemma:

$$\Omega = \{\alpha, \beta\}, \quad \mathcal{B}_1(\alpha) = \{\alpha\}, \mathcal{B}_1(\beta) = \{\beta\}, \quad f_1(\alpha, \{\alpha\}) = f_1(\alpha, \Omega) = \{\alpha\},$$

$$f_1(\beta, \{\beta\}) = f_1(\beta, \Omega) = \{\beta\}, \quad f_1(\alpha, \{\beta\}) = \{\beta\}, \quad f_1(\beta, \{\alpha\}) = \{\alpha\}, \quad \sigma_1(\alpha) = C, \sigma_1(\beta) = D,$$

$$\sigma_2(\alpha) = c, \sigma_2(\beta) = d \text{ (we have omitted the beliefs of Player 2).}$$

At state α Player 1 believes that she is playing C and Player 2 is playing c ($\mathcal{B}_1(\alpha) = \{\alpha\}$ and $\sigma_1(\alpha) = C$, $\sigma_2(\alpha) = c$); furthermore the proposition “Player 1 plays D ” is represented by the event $\{\beta\}$ (this is the only state where Player 1 plays D) and thus, since $f_1(\alpha, \{\beta\}) = \{\beta\}$ and $\sigma_2(\beta) = d$, Player 1 believes that – under the supposition that she plays D – Player 2 plays d and thus her own payoff is 1.

Are the beliefs expressed in COR_2 compatible with rationality? The principles of “rational” belief revision, that are captured by the properties listed in Footnote 9, are principles of logical coherence of dispositional beliefs⁹ and, in general, do not impose any constraints on the content of a counterfactual belief. Thus the above beliefs of Player 1 *could* be rational beliefs, in the sense that they do not violate logical principles or principles of coherence. Those who claim that the beliefs expressed in COR_2 are irrational appeal to the argument that they

⁸ Alternatively, one could remove the initial beliefs $\{\mathcal{B}_i\}_{i \in N}$ from the definition of a model and recover them from the function f_i by taking $f_i(\omega, \Omega)$ to be the set of states that player i – initially – considers possible at state ω . There are further consistency properties that are usually imposed: (1) if $E \neq \emptyset$ then $f_i(\omega, E) \neq \emptyset$, (2) if $\mathcal{B}_i(\omega) \cap E \neq \emptyset$ then $f_i(\omega, E) = \mathcal{B}_i(\omega) \cap E$ and (3) if $E \subseteq F$ and $f_i(\omega, F) \cap E \neq \emptyset$ then $f_i(\omega, E) = f_i(\omega, F) \cap E$. For a more detailed discussion see Bonanno (2013a).

⁹ The principles that were introduced by Alchourrón et al (1985), which pioneered the vast literature on dispositional belief revision, known as the AGM theory.

imply a belief by Player 1 that her “switching” from C to D *causes* Player 2 to change her decision from c to d , while such a causal effect is ruled out by the fact that each player is making her choice in ignorance of the choice made by the other player (the choices are made “simultaneously”). For example, Harper (1988, p. 25) claims that “a causal independence assumption is part of the idealization built into the normal form” and Stalnaker (1996, p. 138) writes “[I]n a strategic form game, the assumption is that the strategies are chosen independently, which means that the choices made by one player cannot influence the beliefs or the actions of the other players”. One can express this point of view as follows (‘IND’ stands for ‘independence’):

In an enriched model of a game, if at state ω player i is rational and chooses strategy x and considers it possible that his opponent is choosing any one of the strategies w_1, \dots, w_n , then the following must be true for every strategy y of player i : (IND _{i})
at any state $\omega' \in f_i(\omega, [y])$ (where $[y]$ denotes the event that –
that is, the set of states where – player i chooses y), player i
continues to consider it possible that his opponent is choosing
any one of the strategies w_1, \dots, w_n and no other strategies.

The beliefs expressed in COR₂ violate the condition IND _{i} .

Should IND _{i} be viewed as a necessary condition for rational beliefs? It seems that, in general, the answer should be negative, for the following reasons.

1. Bicchieri and Green (1999) point out a scenario (an agentive analogue of the Einstein-Podolsky-Rosen phenomenon in quantum mechanics) where causal independence is compatible with correlation and thus it would be possible for a player to coherently believe (a) that her choice is causally independent of the opponent’s choice and also (b) that there is correlation between her choice and the opponent’s choice, such as the correlation expressed in COR₂. A belief of this nature could perhaps be judged to be farfetched or implausible (similarly, perhaps, to the beliefs discussed in the Introduction), but not necessarily irrational.
2. In a series of contributions, Spohn (2003, 2007, 2010, 2012) put forward a new solution concept, which he calls “dependency equilibrium”, which allows for correlation between the players’ choices. An example of a dependency equilibrium is (C, c) (that is,

cooperation) in the Prisoner's Dilemma. Spohn stresses the fact that the notion of dependency equilibrium is consistent with the causal independence of the players' actions:

“The point then is to conceive the decision situations of the players as somehow jointly caused and as entangled in a dependency equilibrium... [B]y no means are the players assumed to believe in a causal loop between their actions; rather, they are assumed to believe in the possible entanglement as providing a common cause of their actions.” (Spohn, 2007, p. 787.)

3. When it comes to judging a player's beliefs about the strategies of her opponents, it is a widely held opinion that it can be fully rational for, say, Player 3 to believe –in a simultaneous game – (a) that the choices of Player 1 and Player 2 are causally independent and yet (b) that “if Player 1 plays x then Player 2 will play x and if Player 1 plays y then Player 2 will play y ”. For example, Aumann (1987, p. 16) writes:

“In a game with more than two players, correlation may express the fact that what 3, say, thinks that 1 will do may depend on what he thinks 2 will do. This has no connection with any overt or even covert collusion between 1 and 2; they may be acting entirely independently. Thus it may be common knowledge that both 1 and 2 went to business school, or perhaps to the same business school; but 3 may not know what is taught there. In that case 3 would think it quite likely that they would take similar actions, without being able to guess what those actions might be.”

Similarly, Brandenburger and Friedenberg (2008, p. 32) write that this correlation in the mind of Player 3 between the action of Player 1 and the action of Player 2 “is really just an adaptation to game theory of the usual idea of common-cause correlation.”

Thus Player 1's beliefs expressed in COR_2 might be criticized for being implausible or hard to justify, but not necessarily irrational.

3.2 Subjective counterfactuals as beliefs about causality

The usual argument in support of the thesis that, for the Prisoner's Dilemma game shown in Figure 1, Player 1's reasoning expressed in COR_1 is fallacious is that even if (e.g. because of symmetry or because of the “identity” assumption) one agrees that the outcome must be one of the two on the diagonal (z_1 and z_4), the off-diagonal outcomes (z_2 and z_3) are nevertheless

causally possible. Thus one must distinguish between *causal* (or objective) possibility and *doxastic* (or subjective) possibility and in the process of rational decision making one has to consider the relevant causal possibilities, even if they are ruled out as doxastically impossible. This is where objective counterfactuals become relevant. This line of reasoning is at the core of causal decision theory.¹⁰

According to this point of view, subjective counterfactuals should be interpreted in terms of the composition of a belief relation \mathcal{B}_i with an objective counterfactual selection function $f: \Omega \times 2^\Omega \rightarrow 2^\Omega$. Under this interpretation, $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} f(\omega', E)$ is the set of states in E that – according to player i 's beliefs at state ω – would be “causally true” if E were the case.

It is worth repeating that – from the point of view of judging the rationality of a choice – what matters is not the “true” causal effect of that choice but what the agent *believes* to be the causal effect of her choice, as illustrated in the example of Section 2 concerning the mislabeled faucets. As another example, consider the case of a player who believes to be engaged – as Player 1 – in a Prisoner's Dilemma game, while in fact Player 2 is a computer that will be informed of Player 1's choice and has been programmed to mirror that choice. In this case, in terms of objective counterfactuals, there is perfect correlation between the choices of the two players, so that the best choice of Player 1 would be to play C . However, Player 1 may rationally play D if she believes that (1) Player 2 will play d and (2) if she were to play C then Player 2 would still play d . Since a player may hold erroneous beliefs about the causal effects of her own choices, it is not clear whether there is a relevant conceptual difference between the “objective” approach discussed in this section and the subjective approach discussed in Section 3.1.¹¹

Causal independence, at a state ω , between the choice of player i and the choices of her opponents would be expressed by the following restriction on the objective selection function

¹⁰ There are various formulations of causal decision theory: Gibbard and Harper (1978), Lewis (1981), Skyrms (1982) and Sobel (1986). For an overview see Weirich (2008).

¹¹ Although in strategic-form games the two approaches may, from a formal point of view, be equivalent, this is not so for dynamic games, where the “objective” approach may be too restrictive. This point is discussed in Bonanno (2013a).

[given a state ω and a player i we denote by $\sigma_{-i}(\omega) = (\sigma_1(\omega), \dots, \sigma_{i-1}(\omega), \sigma_{i+1}(\omega), \dots, \sigma_n(\omega))$ the profile of strategies chosen by the players other than i]:

For every strategy y of player i , if $\omega' \in f(\omega, [y])$ (where $[y]$ denotes the event that – that is, the set of states where – player i chooses y), then $\sigma_{-i}(\omega') = \sigma_{-i}(\omega)$. (IND_2^{obj})

However, as noted above, what matters is not whether IND_2^{obj} holds at state ω but whether player i believes that IND_2^{obj} holds. Hence the following, subjective, version of independence is the relevant condition:

For every strategy y of player i and for every $\omega' \in \mathcal{B}_i(\omega)$, if $\omega'' \in f(\omega', [y])$ then $\sigma_{-i}(\omega'') = \sigma_{-i}(\omega')$. (IND_2^{subj})

It is straightforward to check that condition IND_2^{subj} is equivalent to condition IND_1 if one defines $f_i(\omega, E) = \bigcup_{\omega' \in \mathcal{B}_i(\omega)} f(\omega', E)$, for every event E .

4. Rationality of choice: discussion of the literature

We are yet to provide a precise definition of rationality in strategic-form games. With the few exceptions described below, there has been no formal discussion of the role of counterfactuals in the analysis of strategic-form games. Aumann (1987) was the first to use the notion of epistemic¹² model of a strategic-form game. His definition of rationality, which we will state in terms of beliefs and call Aumann-rationality, is as follows. Recall that, given a state ω in a model of a game and a player i , $\sigma_i(\omega)$ denotes the strategy chosen by player i at state ω and $\sigma_{-i}(\omega) = (\sigma_1(\omega), \dots, \sigma_{i-1}(\omega), \sigma_{i+1}(\omega), \dots, \sigma_n(\omega))$ the profile of strategies chosen by the other players.

¹² The models used by Aumann (1987, 1995) make use of knowledge, that is, of necessarily correct beliefs. We refer to these models as epistemic, reserving the term ‘doxastic’ for models that use the more general notion of belief, which allows for the possibility of error.

Definition 2. Consider a model of a strategic form game (see Definition 1), a state ω and a player i . Player i 's choice at state ω is *Aumann-rational* if there is no other strategy s_i of player i such that $\pi_i(s_i, \sigma_{-i}(\omega')) > \pi_i(\sigma_i(\omega), \sigma_{-i}(\omega'))$ for every $\omega' \in \mathcal{B}_i(\omega)$.¹³ That is, player i 's choice is rational if it is not the case that player i believes that another strategy of hers is strictly better than the chosen strategy.

Note that the above definition is weaker than the definition used in Aumann (1987), since – for simplicity – we have restricted attention to ordinal payoffs and qualitative (that is, non-probabilistic, beliefs).¹⁴ However, *the essential feature of this definition is that it evaluates alternative strategies of player i keeping the beliefs of player i constant*. Hence implicit in this definition of rationality is either a theory of subjective counterfactuals that assumes condition IND_1 or an objective theory of counterfactuals that assumes condition IND_2^{subj} .

The only attempts (that I am aware of) to bring the relevant counterfactuals to the surface are Shin (1992), Bicchieri and Green (1999), Zambrano (2004) and Board (2006).

Shin (1992) develops a framework which is very similar to one based on subjective selection functions (as described in Section 3.1). For each player i in a strategic-form game Shin defines a “subjective state space” Ω_i . A point in this space specifies a belief of player i about his own choice and the choices of the other players. Such belief assigns probability 1 to player i 's own choice (that is, player i is assumed to know his own choice). Shin then defines a metric on this space as follows. Let ω be a state where player i attaches probability 1 to his own choice, call it A , and has beliefs represented by a probability distribution P on the strategies of his opponents; the closest state to ω where player i chooses a different strategy, say B , is a state ω'

¹³ Recall that $\mathcal{B}_i(\omega)$ is the set of states that player i considers possible at state ω ; recall also the assumption that $\sigma_i(\bullet)$ is constant on $\mathcal{B}_i(\omega)$, that is, for every $\omega' \in \mathcal{B}_i(\omega)$, $\sigma_i(\omega') = \sigma_i(\omega)$.

¹⁴ When payoffs are taken to be von Neumann-Morgenstern payoffs and the beliefs of player i at state ω are represented by a probability distribution $\mathbf{p}_{i,\omega} : \Omega \rightarrow [0,1]$ (assuming that Ω is a finite set) whose support coincides with $\mathcal{B}_i(\omega)$ (that is, $\mathbf{p}_{i,\omega}(\omega') > 0$ if and only if $\omega' \in \mathcal{B}_i(\omega)$) then the choice of player i at state ω is

where player i attaches probability 1 to B and has the same probability distribution P over the strategies of his opponents that he has at ω . This metric allows player i to evaluate the counterfactual “if I chose B then my payoff would be x ”. Thus Shin imposes *as an axiom* the requirement that player i should hold the same beliefs about the other players’ choices when contemplating a “deviation” from his actual choice. This assumption corresponds to requirement IND_1 . Not surprisingly, his main result is that a player is rational with respect to this metric if and only if she is Aumann-rational.

Zambrano’s (2004) approach is a mixture of objective and subjective counterfactuals. His analysis is restricted to two-player strategic-form games. First of all, he defines a “subjective” selection function for player i , $f_i: \Omega \times S_i \rightarrow \Omega$, which follows Stalnaker (1968) in assuming that, for every hypothesis and every state ω , there is a *unique* world closest to ω where that hypothesis is satisfied; furthermore, the hypotheses consist of the possible strategies of player i (the set of strategies S_i), rather than events. He interprets $f_i(\omega, s_i) = \omega'$ as follows: “state ω' is the state closest to ω , according to player i , in which player i deviates from the strategy prescribed by ω and, instead, plays s_i ” (p. 5). He then imposes the requirement that “player i is the *only* one that deviates from $\sigma(\omega)$ in $f_i(\omega, s_i)$, that is, $\sigma_j(f_i(\omega, s_i)) = \sigma_j(\omega)$ ” (Condition F2, p. 5; j denotes the other player). This appears to be in the spirit of the objective causal independence assumption IND_2^{obj} . However, Zambrano does not make use of this requirement, because he focuses on the beliefs of player i at the state $f_i(\omega, s_i)$ and uses these beliefs to evaluate *both* the original strategy $\sigma_i(\omega)$ *and* the new strategy s_i . He introduces the following definition of rationality:

“player i is W-rational [at state ω] if there is no deviation $s_i \neq \sigma_i(\omega)$ such that strategy s_i is preferred to $\sigma_i(\omega)$ given the belief that player i holds *at the state closest to ω in which i deviates to s_i* . The interpretation is that the rationality of choosing strategy $\sigma_i(\omega)$ at state ω against a deviation $s_i \neq \sigma_i(\omega)$ is determined with respect to beliefs

defined to be rational if and only if it maximizes player i ’s expected payoff at state ω , that is, if and only if there is no strategy s_i of player i such that
$$\sum_{\omega' \in \mathcal{B}_i(\omega)} \mathbf{p}_{i,\omega}(\omega') \pi_i(s_i, \sigma_i(\omega')) > \sum_{\omega' \in \mathcal{B}_i(\omega)} \mathbf{p}_{i,\omega}(\omega') \pi_i(\sigma_i(\omega), \sigma_i(\omega')).$$

that arise at the closest state to ω in which s_i is actually chosen, that is, with respect to beliefs at $f_i(\omega, s_i)$.” (Zambrano, 2004, p. 6).

Expressed in terms of our qualitative approach, player i is W-rational at state ω if there is no strategy s_i of player i such that $\pi_i(s_i, \sigma_{-i}(\omega')) > \pi_i(\sigma_i(\omega), \sigma_{-i}(\omega'))$ for every $\omega' \in \mathcal{B}_i(f_i(\omega, s_i))$. Note that, unlike Aumann-rationality (Definition 2), the quantification is over $\mathcal{B}_i(f_i(\omega, s_i))$ rather than over $\mathcal{B}_i(\omega)$.¹⁵ The definition of W-rationality thus disregards the beliefs of player i at state ω and focuses instead on the beliefs that player i would have if she changed her strategy. Since, in general, those hypothetical beliefs can be different from the initial beliefs at state ω , there is no connection between W-rationality and Aumann-rationality. For example, consider the game shown in Part *a* of Figure 3 and the model shown in Part *b*.

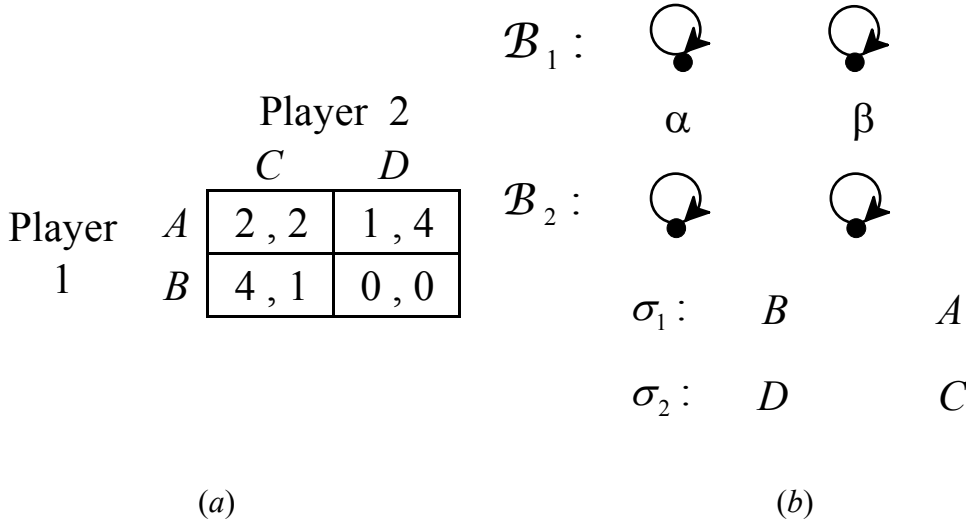


Figure 2
(a) A strategic-form game. (b) A model of the game.

¹⁵ Zambrano uses probabilistic beliefs: for every $\omega \in \Omega$, $\mathbf{p}_{i,\omega} : \Omega \rightarrow [0,1]$ is a probability distribution over Ω that represents the beliefs of player i at state ω . Our set $\mathcal{B}_i(\omega)$ corresponds to the support of $\mathbf{p}_{i,\omega}$.

Zambrano's definition is as follows: player i is W-rational at state ω if there is no strategy s_i of player i such that

$$\sum_{\omega' \in \Omega} \mathbf{p}_{i,f_i(\omega,s_i)}(\omega') \pi_i(s_i, \sigma_j(\omega')) > \sum_{\omega' \in \Omega} \mathbf{p}_{i,f_i(\omega,s_i)}(\omega') \pi_i(\sigma_i(\omega), \sigma_j(\omega')).$$

Let the selection function of Player 1 be given by $f_1(\alpha, B) = f_1(\beta, B) = \alpha$ and $f_1(\alpha, A) = f_1(\beta, A) = \beta$. Consider state α where the play is (B, D) and both players get a payoff of 0. Player 1 is W-rational at state α (where she chooses B and believes that Player 2 chooses D) because if she were to play A (state β) then she would believe that Player 2 played C and – given these beliefs – playing B is better than playing A . However, Player 1 is not Aumann-rational at state α , because the notion of Aumann rationality uses the beliefs of Player 1 at state α to compare A to B (while the notion of W-rationality uses the beliefs at state β).

Zambrano then shows (indirectly, through the implications of common knowledge of rationality) that W-rationality coincides with Aumann-rationality if one adds the following restriction to the subjective selection function f_i : for every state ω and every strategy $s_i \in S_i$, $\text{marg}_{S_j} p_{i,\omega}(\bullet) = \text{marg}_{S_j} p_{i,f_i(\omega, s_i)}(\bullet)$, that is, at the closest state to ω where player i plays strategy s_i , the beliefs of player i concerning the strategy chosen by the other player (player j) are the same as at state ω . This is, of course, condition IND_1 .

Board (2006) uses objective counterfactuals as defined by Stalnaker (1968) (for every hypothesis and every state ω , there is a *unique* world closest to ω where that hypothesis is satisfied). Like Zambrano, Board takes as possible hypotheses the individual strategies of the players: he introduces an objective selection function $f: \Omega \times \bigcup_{i \in N} S_i \rightarrow \Omega$, that specifies – for every state ω , every player i and every strategy $s_i \in S_i$ of player i – the unique world $f(\omega, s_i) \in \Omega$ closest to ω where player i chooses s_i . Recall that $\sigma_i(\omega)$ denotes the strategy chosen by player i at state ω . In accordance with Stalnaker's theory of counterfactuals, Board assumes that $f(\omega, \sigma_i(\omega)) = \omega$, that is, the closest state to ω where player i chooses the strategy that he chooses at ω is ω itself. On the other hand, if $s_i \neq \sigma_i(\omega)$ and $f(\omega, s_i) = \omega'$ then it is necessarily the case that $\omega' \neq \omega$, since it must be that $\sigma_i(\omega') = s_i$. What does player i believe at state ω about the choices of the other players? As before, let \mathcal{B}_i be the belief relation of player i and $\mathcal{B}_i(\omega) = \{\omega' \in \Omega: \omega \mathcal{B}_i \omega'\}$ the belief set of player i at state ω . We denote by $S_{-i} = S_1 \times \dots \times S_{i-1} \times S_{i+1} \times \dots \times S_n$ the set of strategy profiles for the players other than i . Then the strategy profiles of the opponents that player i considers possible at state ω , if she plays her

chosen strategy $\sigma_i(\omega)$, is $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(\omega')\} = \{s_{-i} \in S_{-i} : s_{-i} = \sigma_{-i}(\omega') \text{ for some } \omega' \in \mathcal{B}_i(\omega)\}$. On

the other hand, what are her beliefs – at state ω – about the strategy profiles of her opponents if she were to choose a strategy $s_i \neq \sigma_i(\omega)$? For every state ω' that she deems possible at state ω (that is, for every $\omega' \in \mathcal{B}_i(\omega)$) she considers the closest state to ω' where she plays s_i , namely $f(\omega', s_i)$, and looks at the choices made by her opponents at state $f(\omega', s_i)$.¹⁶ Thus the strategy profiles of the opponents that player i would consider possible at state ω , if she were to play a strategy $s_i \neq \sigma_i(\omega)$, is $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(f(\omega', s_i))\}$. Note that, in general, there is no relationship between the sets $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(f(\omega', s_i))\}$ and $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(\omega')\}$; indeed, these two sets might even be disjoint.

Board defines player i to be *causally rational* at state ω (where she chooses strategy $\sigma_i(\omega)$) if it is not the case that she believes, at state ω , that there is another strategy $s_i \in S_i$ which would yield a higher payoff than $\sigma_i(\omega)$. His definition is expressed in terms of expected payoff maximization.¹⁷ Since, in general, the two sets $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(f(\omega', s_i))\}$ and

¹⁶ Recall the assumption that a player always knows her chosen strategy, that is, for every $\omega' \in \mathcal{B}_i(\omega)$, $\sigma_i(\omega') = \sigma_i(\omega)$ and thus – since we are considering a strategy $s_i \neq \sigma_i(\omega)$ – it must be the case that $f(\omega', s_i) \neq \omega'$.

¹⁷ Like Zambrano, Board assumes that payoffs are von Neumann-Morgenstern payoffs and beliefs are probabilistic: for every $\omega \in \Omega$, $\mathbf{p}_{i,\omega}$ is a probability distribution over $\mathcal{B}_i(\omega)$ that represents the probabilistic beliefs of player i at state ω . Board defines player i to be causally rational at state ω if there is no strategy s_i that would yield a higher expected payoff if chosen instead of $\sigma_i(\omega)$, that is, if there is no $s_i \in S_i$ such that
$$\sum_{\omega' \in \mathcal{B}_i(\omega)} \mathbf{p}_{i,\omega}(\omega') \pi_i(s_i, \sigma_{-i}(f(\omega', s_i))) > \sum_{\omega' \in \mathcal{B}_i(\omega)} \mathbf{p}_{i,\omega}(\omega') \pi_i(\sigma_i(\omega), \sigma_{-i}(\omega'))$$
. There is no clear qualitative counterpart to this definition, because of the lack of any constraints that relate $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(f(\omega', s_i))\}$ to $\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(\omega')\}$. Board (2006, p. 16) makes this point as follows: “since each state describes what each player does as well as what her opponents do, the player will change the state if she changes her choice. There is no guarantee that her opponents will do the same in the new state as they did in the original state.”

$\bigcup_{\omega' \in \mathcal{B}_i(\omega)} \{\sigma_{-i}(\omega')\}$ might be disjoint, causal rationality is consistent with each player choosing cooperation in the Prisoner's Dilemma. To see this, consider the following model of the Prisoner's Dilemma game of Figure 1: $\Omega = \{\alpha, \beta\}$, $\sigma_1(\alpha) = C$, $\sigma_1(\beta) = D$, $\sigma_2(\alpha) = c$, $\sigma_2(\beta) = d$, $\mathcal{B}_1(\alpha) = \mathcal{B}_2(\alpha) = \{\alpha\}$, $\mathcal{B}_1(\beta) = \mathcal{B}_2(\beta) = \{\beta\}$, $f(\alpha, C) = f(\alpha, c) = f(\beta, C) = f(\beta, c) = \alpha$ and $f(\beta, D) = f(\beta, d) = f(\alpha, D) = f(\alpha, d) = \beta$. Then at state α Player 1 is causally rational: she chooses C and believes that her payoff will be 2 (because she believes that Player 2 has chosen c) and she also believes that if she were to play D then Player 2 would play d and thus her payoff would be 1. Note that this model is a formal representation of the reasoning expressed in COR₁. Board's main result is that a necessary and sufficient condition for causal rationality to coincide with Aumann rationality is the IND_2^{subj} condition of Section 3.2.¹⁸

Bicchieri and Green's (1999) aim is to clarify the implications of the "Identity assumption" in the Prisoner's Dilemma game. They enrich the definition of a model of a game (Definition 1) by adding a binary relation $C \subseteq \Omega \times \Omega$ of "nomic accessibility", interpreting $\omega C \omega'$ as " ω' is causally possible relative to ω " in the sense that "everything that occurs at ω' is consistent with the laws of nature that hold at ω " (p. 180). After discussing at length the difference between doxastic possibility (represented by the relations \mathcal{B}_i , $i \in N$) and causal possibility (in the spirit of causal decision theory), they raise the question whether it is possible to construe a situation in which it is causally necessary that the choices of the two players in the Prisoner's Dilemma are the same, while their actions are nonetheless causally independent. They suggest that the answer is positive: one could construct an agentive analogue of the Einstein-Podolsky-Rosen phenomenon in quantum mechanics (p. 184). They conclude that there may indeed be a coherent nomic interpretation of the Identity assumption, but such interpretation may be controversial.

¹⁸ Board presents this as an objective condition on the selection function (if $\omega' = f(\omega, s_i)$ then $\sigma_{-i}(\omega') = \sigma_{-i}(\omega)$) assumed to hold at every state but then acknowledges (p. 12) that "it is players' beliefs in causal independence rather than causal independence itself that drives the result."

In the next section we re-examine the commonly held view that counterfactuals ought to be considered inherent to decision-making and rationality.

5. Conditionals of deliberation and pre-choice beliefs

A common feature of all the epistemic/doxastic models of games used in the literature is the assumption that if a player takes a particular action at state ω then she knows, at state ω , that she takes that action. This approach thus requires the use of either objective or subjective counterfactuals in order to represent a player's beliefs about the consequences of taking alternative actions. However, several authors have pointed out that *it is the essence of deliberation that one cannot reason towards a choice if one already knows what that choice will be*. For instance, Shackle (1958, p. 21) remarks that if an agent could predict the option he will choose, his decision problem would be “empty”, Ginét (1962, p. 50) claims that “it is conceptually impossible for a person to know what a decision of his is going to be before he makes it”, Goldman (1970, p. 194) writes that “deliberation implies some doubt as to whether the act will be done”, Levi states that “the deliberating agent cannot, before choice, predict how he will choose” (Levi, 1986, p. 65) and coins the phrase “deliberation crowds out prediction” (Levi, 1997, p. 81), Spohn (2012, p. 109) writes that “the decision model must not impute to the agent any cognitive or doxastic assessment of his own actions”.¹⁹

Deliberation involves reasoning along the following lines: “if I take action a , then the outcome will be x and if I take action b , then the outcome will be y ”. Indeed it has been argued (DeRose, 2010) that the appropriate conditionals for deliberation are *indicative* conditionals, rather than subjunctive conditional. If I say “if I had left the office at 4 pm I would not have been stuck in traffic”, I convey the information that – as a matter of fact – I did not leave the office at 4 pm and thus I am uttering a counterfactual conditional, namely one which has a false antecedent (such a statement would not make sense if uttered before 4 pm). On the other hand, if I say “if I leave the office at 4 pm I will not be stuck in traffic” I am uttering what is normally called an indicative conditional and I am conveying the information that I am evaluating the

¹⁹ Similar observations can be found in Gilboa (1999), Schick (1979), Spohn (1977), Kadane and Seidenfeld (1999). For a discussion and further references see Ledwig (2005).

consequences of a possible future action (such a statement would not make sense if uttered after 4 pm). Concerning the latter conditional, is there a difference between the indicative mood and the subjunctive mood? If I said (before 4 pm) “if I were to leave the office at 4 pm I would not be stuck in traffic”, would I be conveying the same information as with the previous indicative conditional? On this point there does not seem to be a clear consensus in the literature. I agree with DeRose's claim that the subjunctive mood conveys different information relative to the indicative mood: its role is to

“call attention to the possibility that the antecedent is (or will be) false, where one reason one might have for calling attention to the possibility that the antecedent is (or will be) false is that it is quite likely that it is (or will be) false.” (DeRose, 2010, p. 10.)

The indicative conditional signals that the decision whether to leave the office at 4 pm is still “open”, while the subjunctive conditional intimates that the speaker is somehow ruling out that option: for example, he has made a tentative or firm decision not to leave at 4 pm.

Thus it would be desirable to model a player's *deliberation* (or *pre-choice*) stage beliefs, where the player considers the consequences of all her actions, *without predicting her subsequent decision*. If a state encodes the player's actual choice, then that choice can be judged to be rational or irrational by relating it to the player's pre-choice beliefs. Hence, according to this approach, it becomes possible for a player to have the same beliefs in two different states, ω and ω' , and be labeled as rational at state ω and irrational at state ω' , because the action she ends up taking at state ω is optimal given those beliefs, while the action she ends up taking at state ω' is not optimal given those same beliefs.

A potential objection to this view arises in dynamic games where a player chooses more than once along a given play of the game. Consider a situation where at time t_1 player i faces a choice and knows that she might be called upon to make a second choice at a later time t_2 . The view outlined above requires player i to have “open” beliefs about her choice at time t_1 but also allows her to have beliefs about (or be certain of) what choice she will make at the later time t_2 . Is this problematic? Several authors have maintained that there is no inconsistency between the principle that one should not attribute to a player beliefs about her current choice and the claim

that, on the other hand, one can attribute to the player beliefs about her later choices. For example, Gilboa writes:

“[W]e are generally happier with a model in which one cannot be said to have beliefs about (let alone knowledge of) one’s own choice while making this choice . [O]ne may legitimately ask: Can you truly claim you have no beliefs about your own future choices? Can you honestly contend you do not believe – or even know – that you will not choose to jump out of the window? [T]he answer to these questions is probably a resounding “No”. But the emphasis should be on timing: when one considers one’s choice tomorrow, one may indeed be quite sure that one will not decide to jump out of the window. However, a future decision should actually be viewed as a decision by a different “agent” of the same decision maker. [...] It is only at the time of choice, within an “atom of decision”, that we wish to preclude beliefs about it.” (Gilboa,1999, pp. 171 –172)

In a similar vein, Levi (1997 , p. 81) writes that “agent X may coherently assign unconditional credal probabilities to hypotheses as to what he will do when some future opportunity for choice arises. Such probability judgments can have no meaningful role, however, when the opportunity of choice becomes the current one.” Similarly, Spohn (1977, p. 114) states the principle that “any adequate quantitative decision model must not explicitly or implicitly contain any subjective probabilities for acts” and then adds (Spohn, 1999, pp. 44 –45) that in the case of sequential decision making, the decision maker can ascribe subjective probabilities to his future (but not to his present) actions. We share the point of view expressed by these authors. If a player moves sequentially at times t_1 and t_2 , with $t_1 < t_2$, then at time t_1 she has full control over her immediate choices (those available at t_1) but not over her later choices (those available at t_2). The agent can predict – or form an intention about – her future behavior, but she cannot irrevocably decide it, just as she can predict – but not decide – how other individuals will behave after her current choice.

Doxastic models of games incorporating deliberation-stage beliefs were recently introduced in Bonanno (2013*b*, 2013*c*) for the analysis of dynamic games. Space limitations prevent us from going into the details of these models.

References

- Alchourrón, Carlos, Gärdenfors, Peter and Makinson, David 1985. "On the logic of theory change: partial meet contraction and revision functions", *The Journal of Symbolic Logic*, 50:510-530.
- Arló-Costa, Horacio and Bicchieri, Cristina 2007. "Knowing and supposing in games of perfect information", *Studia Logica*, 86: 353-373.
- Aumann, Robert 1987. "Correlated equilibrium as an expression of Bayesian rationality", *Econometrica*, 55: 1-19.
- Aumann, Robert 1995. "Backward induction and common knowledge of rationality," *Games and Economic Behavior*, 8: 6-19.
- Battigalli, Pierpaolo and Bonanno, Giacomo 1999. "Recent results on belief, knowledge and the epistemic foundations of game theory", *Research in Economics*, 53:149-225.
- Bicchieri, Cristina and Green, Mitchell 1999. "Symmetry arguments for cooperation in the Prisoner's Dilemma" in C. Bicchieri, R. Jeffrey and B. Skyrms (eds.), *The logic of strategy*, Oxford University Press, pp. 175-195.
- Binmore, Ken 2011. *Rational decisions*, Princeton University Press.
- Board, Oliver 2004. "Dynamic interactive epistemology", *Games and Economic Behavior*, 49: 49-80.
- Board, Oliver 2006. "The equivalence of Bayes and causal rationality in games", *Theory and Decision*, 61: 1-19.
- Bonanno, Giacomo 2013a. "Reasoning about strategies and rational play in dynamic games" in J. van Benthem, S. Ghosh and R. Verbrugge (eds.), *Modeling strategic reasoning*, Texts in Logic and Games, Springer, forthcoming.
- Bonanno, Giacomo 2013b. "A dynamic epistemic characterization of backward induction without counterfactuals", *Games and Economic Behavior*, 78: 31-43.
- Bonanno, Giacomo 2013c. "An epistemic characterization of generalized backward induction", Working Paper No. 134, University of California Davis (<http://ideas.repec.org/p/cda/wpaper/13-4.html>).
- Brams, Steven 1975. "Newcomb's Problem and Prisoners' Dilemma", *The Journal of Conflict Resolution*, 19:496-612.
- Brandenburger, Adam and Friedenberg, Amanda 2008. "Intrinsic correlation in games", *Journal of Economic Theory*, 141: 28-67.

- Clausing, Thorsten 2004. "Belief revision in games of perfect information", *Economics and Philosophy*, 20: 89-115.
- Davis, Lawrence 1977. "Prisoners, paradox and rationality", *American Philosophical Quarterly*, 14: 319-327.
- Davis, Lawrence 1985. "Is the symmetry argument valid?" in R. Campbell and L. Snowden (eds.), *Paradoxes of rationality and cooperation*, University of British Columbia Press, pp. 255-262.
- DeRose, Keith 2010. "The conditionals of deliberation", *Mind*, 119: 1-42.
- Gibbard, Allan and Harper, William 1978. "Counterfactuals and two kinds of expected utility" in W. Harper, R. Stalnaker and G. Pearce (eds.), *Ifs: conditionals, belief, decision, chance, and time*, D. Reidel, pp. 153-190.
- Gilboa, Itzhak 1999. "Can free choice be known?" in C. Bicchieri, R. Jeffrey and B. Skyrms (eds.), *The logic of strategy*, Oxford University Press, pp. 163-174.
- Ginet, Carl 1962. "Can the will be caused?", *The Philosophical Review*, 71: 49-55.
- Goldman, Alvin 1970. *A theory of human action*, Princeton University Press
- Halpern, Joseph 1999. "Hypothetical knowledge and counterfactual reasoning", *International Journal of Game Theory*, 28: 315-330.
- Halpern, Joseph 2001. "Substantive rationality and backward induction", *Games and Economic Behavior*, 37: 425-435.
- Harper, William L. 1988. "Causal decision theory and game theory: a classic argument for equilibrium solutions, a defense of weak equilibria, and a limitation for the normal form representation" in W. L. Harper and B. Skyrms (eds.), *Causation in decision, belief change, and statistics*, II, Kluwer Academic Publishers, pp. 246-266.
- Hausman, Daniel 2012. *Preference, value, choice and welfare*, Cambridge University Press.
- Kadane, Joseph B. and Seidenfeld, Teddy 1999. "Equilibrium, common knowledge, and optimal sequential decisions" in J. B. Kadane, Mark J. Schervish and T. Seidenfeld (eds.), *Rethinking the Foundations of Statistics*, Cambridge University Press, pp. 27-46.
- Ledwig, Marion 2005. "The no probabilities for acts-principle", *Synthese*, 144: 171--180
- Levi, Isaac 1986. *Hard choices*, Cambridge University Press.
- Levi, Isaac 1997. *The covenant of reason: rationality and the commitments of thought*, Cambridge University Press.
- Lewis, David 1973. *Counterfactuals*, Oxford, Basil Blackwell.
- Lewis, David 1981. "Causal decision theory", *Australasian Journal of Philosophy*, 59: 5-30.

- Rabinowicz, Wlodek 2000. "Backward induction in games: on an attempt at logical reconstruction" in W. Rabinowicz (ed.), *Value and choice: some common themes in decision theory and moral philosophy*, University of Lund Philosophy Reports, pp. 243-256.
- Rubinstein, Ariel 2012. *Economic fables*, Open Book Publishers.
- Rubinstein, Ariel and Salant, Yuval 2008. "Some thoughts on the principle of revealed preference" in A. Caplin and A. Schotter (eds.), *Handbook of economic methodology*, Oxford University Press, pp. 116-124.
- Shackle, George L.S. 1958. *Time in Economics*, North-Holland Publishing Company, Amsterdam.
- Shick, Frederic 1979. "Self knowledge, uncertainty and choice", *British Journal for the Philosophy of Science*, 30: 235-252.
- Shin, Hyun Song 1992. "Counterfactuals and a theory of equilibrium in games" in C. Bicchieri and M. L. Dalla Chiara (eds.), *Knowledge, belief and strategic interaction*, Cambridge University Press, pp. 397-413.
- Skyrms, Brian 1982. "Causal decision theory", *Journal of Philosophy*, 79: 695-711.
- Sobel, Jordan H. 1986. "Notes on decision theory: old wine in new bottles", *Australasian Journal of Philosophy*, 64: 407-437.
- Spohn, Wolfgang 1977. "Where Luce and Krantz do really generalize Savage's decision model", *Erkenntnis*, 11: 113-134.
- Spohn, Wolfgang 1999. *Strategic Rationality*, Volume 24 of *Forschungsberichte der DFG-Forschergruppe Logik in der Philosophie*, Konstanz University.
- Spohn, Wolfgang 2003. "Dependency equilibria and the causal structure of decision and game situations", *Homo Oeconomicus*, 20: 195-255.
- Spohn, Wolfgang 2007. "Dependency equilibria", *Philosophy of Science*, 74: 775-789.
- Spohn, Wolfgang 2010. "From Nash to dependency equilibria" in G. Bonanno, B. Löwe and W. van der Hoek (eds.), *Logic and the foundations of game and decision theory – LOFT8*, Texts in Logic and Games, Springer, pp. 135-150.
- Spohn, Wolfgang 2012. "Reversing 30 years of discussion: why causal decision theorists should one-box", *Synthese*, 187: 95-122.
- Stalnaker, Robert 1968. "A theory of conditionals" in N. Rescher (ed.), *Studies in logical theory*, Oxford, Blackwell, pp. 98-112.
- Stalnaker, Robert 1996. "Knowledge, belief and counterfactual reasoning in games", *Economics and Philosophy*, 12: 133-163.

- Weirich, Paul 2008. “Causal decision theory” in *Stanford Encyclopedia of Philosophy*, URL: <http://plato.stanford.edu/entries/decision-causal/>.
- Zambrano, Eduardo 2004. “Counterfactual reasoning and common knowledge of rationality in normal form games”, *Topics in Theoretical Economics*, 4 (1), article 8.