

Dependency Equilibria

Wolfgang Spohn[†]

This paper introduces a new equilibrium concept for normal form games called dependency equilibrium; it is defined, exemplified, and compared with Nash and correlated equilibria in Sections 2–4. Its philosophical motive is to rationalize cooperation in the one shot prisoners’ dilemma. A brief discussion of its meaningfulness in Section 5 concludes the paper.

1. Introduction. In this note I would like to present and briefly discuss a new equilibrium concept for game theory that I call dependency equilibrium. When it occurred to me 24 years ago, I put it aside because it seemed to me of doubtful sense. I do not know whether anybody had the same idea; if so, he or she may have dismissed it for the same reason. In the meantime, I have changed my mind; I think it can be backed up by a meaningful story. Hence, I think the concept at least deserves a hearing, even though the longer story can at best be feebly indicated here.

The driving force behind this concept is, once more, the great riddle posed by the Prisoners’ Dilemma (PD). This has elicited a vast literature and a large number of astonishingly varied attempts to undermine defection as the only rational solution and establish cooperation as a rational possibility, at least in the iterated case. But the hard case, it seems to me, still stands unshaken. Under appropriate conditions backward induction is valid;¹ hence, given full rationality (instead of some form of ‘bounded rationality’) and sufficient common knowledge, continued defection is the only solution in the finitely iterated PD. The same conclusion is reached via the iterated elimination of weakly dominated strategies.² I find this

[†]To contact the author, please write to: Department of Philosophy, University of Konstanz, 78457 Konstanz, Germany; e-mail: Wolfgang.Spohn@uni-konstanz.de.

1. Cf. Aumann 1995, but see also the excellent discussion in Rabinowicz 2000.

2. Iterated elimination of weakly dominated strategies is a reasonable procedure when applied to the iterated PD, all the more so as the criticisms this may meet compared to the elimination of strongly dominated strategies do not obtain in this application. See, e.g., Myerson 1991, Sections 2.5 and 3.1.

Philosophy of Science, 74 (December 2007) pp. 775–000. 0031-8248/2007/7405-0019\$10.00
Copyright 2007 by the Philosophy of Science Association. All rights reserved.

conclusion highly disconcerting; it amounts to an outright refutation of the underlying theory of rationality. Moreover, I find that all the sophisticated observations made so far about PD have failed to tone down this harsh conclusion. Cooperation *must* remain at least a rational possibility in the finitely iterated PD, and under ideal conditions *even more so* than under less ideal ones. Thus, something needs to be changed in standard rationality theory, that is, decision and game theory. After a long time of thinking otherwise, I have come to the conclusion that it is the one-shot case that needs to be reconsidered, and this is what I try to do with the new equilibrium concept. Sections 2–4 will introduce and exemplify the new concept and offer a few technical observations. Section 5 is devoted to a brief discussion.

2. Nash, Correlated, and Dependency Equilibria. Here is an outline of the new concept. For comparison, it is useful to rehearse Nash equilibria and Aumann's correlated equilibria. We shall deal only with normal form games. Hence, the refinements of Nash equilibria relating to the extensive form are outside of our focus. It suffices to consider two-person games. While I hardly develop the theory here, it may be routinely extended, it seems, to n -person games.

Thus, let $A = \{a_1, \dots, a_m\}$ be the set of pure strategies of Ann (row chooser) and $B = \{b_1, \dots, b_n\}$ the set of pure strategies of Bob (column chooser). Let u and v be the utility functions of Ann and Bob, respectively, from $A \times B$ into \mathbb{R} ; we abbreviate $u(a_i, b_k) = u_{ik}$ and $v(a_i, b_k) = v_{ik}$.

Moreover, let S be the set of mixed strategies of Ann, that is, the set of probability distributions over A . Hence, $s = \langle s_1, \dots, s_m \rangle = (s_i) \in S$ if and only if $s_i \geq 0$ for $i = 1, \dots, m$ and $\sum_{i=1}^m s_i = 1$. Likewise, let T be the set of mixed strategies of Bob. Mixed strategies have an ambiguous interpretation. Usually, the probabilities are thought to be intentional mixtures by each player. But it is equally appropriate to interpret them as representing the beliefs of others about the player. Indeed, in relation to dependency equilibria, this will be the only meaningful interpretation.

We shall envisage the possibility that the actions in a game may be governed by any probability distribution whatsoever. Let P be the set of distributions over $A \times B$. Thus, $p = (p_{ik}) \in P$ if and only if $p_{ik} \geq 0$ for all $i = 1, \dots, m$ and $k = 1, \dots, n$ and $\sum_{i,k} p_{ik} = 1$. Each $p \in P$ has a marginal s over A and a marginal t over B . But since p may contain arbitrary dependencies between A and B , it is usually not the product of the marginals s and t . This is all the terminology we shall need.

As is well known, $\langle s, t \rangle \in S \times T$ is defined as a *Nash equilibrium* if and only if for all $j = 1, \dots, m$, $\sum_{i,k} s_i t_k u_{ik} \geq \sum_k t_k u_{jk}$ (or, equivalently, for all $s^* \in S$, $\sum_{i,k} s_i t_k u_{ik} \geq \sum_{i,k} s_i^* t_k u_{ik}$), and if the corresponding condition holds for the other player. Hence, in a Nash equilibrium, neither Ann nor Bob

can raise her or his expected utility by changing from her or his equilibrium strategy to some other pure or mixed strategy, given that the other player sticks to his or her equilibrium strategy. There is no need here to rehearse the standard rationale for Nash equilibria, and there is no time to discuss their strengths and weaknesses.³

Obviously Ann's and Bob's choices from A and B are independent in a Nash equilibrium. This is an assumption I would like to abandon (for reasons that will become clear later on). Aumann (1974) has introduced an equilibrium concept that allows for dependence between the players. Here is the definition from Aumann 1987 (which is a bit simpler and less general than his original definition, the statement of which would require us to introduce additional structure): Let $p \in P$ have marginals $s \in S$ and $t \in T$. Then p is a *correlated equilibrium* if and only if for all $j = 1, \dots, m$, $\sum_{i,k} p_{ik} u_{ik} \geq \sum_k t_k u_{jk}$ (or, equivalently, for all $s^* \in S$, $\sum_{i,k} p_{ik} u_{ik} \geq \sum_{i,k} s_i^* t_k u_{ik}$) and if the corresponding condition holds for the other player. The most straightforward way to understand this, which is offered by Aumann himself (1987, 3ff.), is the following: Somehow, Ann and Bob agree on a joint distribution over the strategy combinations or outcomes of their game. One combination is chosen at random according to this distribution, and each player is told only his or her part of the combination. If no player can raise his or her expected utility by breaking his or her part of the agreed joint distribution and choosing some other pure or mixed strategy instead, then this joint distribution is a correlated equilibrium. Thus, correlated equilibria are self enforcing, they do not need external help from sanctions or agreements.

Correlated equilibria appear to fall outside noncooperative game theory. However, one can model the selection of a joint distribution for the original game as an additional move in a game enlarged by preplay communication, and it then turns out that all and only the Nash equilibria of the enlarged game correspond to correlated equilibria in the original game.⁴ This reflects the fact that correlated equilibria, despite their allowance of dependence, are still noncooperative in essence. The players' standard of comparison is still whether they might be better off by independently doing something else, where their expectations about the other player are given by the marginal over their strategies.

This standard of comparison is changed in the dependency equilibria introduced below. It is not the expected utility given the marginal for the other player, but rather the *conditional expected utility* given the conditional probabilities determined by the joint distribution.

3. This has been done many times, also by myself in Spohn 1982.

4. For details, see Myerson 1991, 255–257.

Here is a first attempt to formalize this idea: Let $p \in P$ have marginals $s \in S$ and $t \in T$. Let $p_{k|i}$ be the probability of b_k given a_i (i.e., $p_{k|i} = p_{ik}/s_i$) and $p_{i|k} = p_{ik}/t_k$ the probability for a_i given b_k . Now, p is a *dependency equilibrium* if and only if for all i with $s_i > 0$ and all $j = 1, \dots, m$, $\sum_k p_{k|i} u_{ik} \geq \sum_k p_{k|j} u_{jk}$ and if the corresponding condition holds for the other player. Thus, in a dependency equilibrium each player maximizes their conditional expected utility with whatever they do with positive probability according to the joint equilibrium distribution.

This provokes at least three immediate remarks. The first point to be taken up is a technical flaw in the above definition. If some a_j has probability 0 in the joint distribution p , that is, if $s_j = 0$, then no conditional probability given a_j is defined. Yet, the fact that $s_j = 0$ should not render the other figures meaningless. This kind of problem is standardly solved by engaging in epsilontics, that is, by approaching probability 0 by ever smaller positive probabilities. This strategy is easily applied here. Let us call a distribution $p \in P$ *strictly positive* if and only if $p_{ik} > 0$ for all i and k . Now we correct my flawed definition by an approximating sequence of strictly positive distributions; this is my official definition: $p \in P$ is a *dependency equilibrium* if and only if there is a sequence $(p^r)_{r \in \mathbb{N}}$ of strictly positive distributions such that $\lim_{r \rightarrow \infty} p^r = p$ and for all i with $s_i > 0$ and $j = 1, \dots, m$, $\lim \sum_k p_{k|i}^r u_{ik} \geq \lim \sum_k p_{k|j}^r u_{jk}$ and for all k with $t_k > 0$ and all $l = 1, \dots, n$, $\lim \sum_i p_{i|k}^r v_{ik} \geq \lim \sum_i p_{i|l}^r v_{il}$. All the conditional probabilities appearing in this definition are well defined. Though the definition looks more complicated now, the intuitive characterization given above still fits perfectly.

After this correction, the second point is that dependency equilibria seem to be well in line with decision theory. Most textbooks state that the general decision rule is to maximize *conditional* expected utility. Savage (1954) still assumed a clear separation of states of the world having probabilities and consequences carrying utilities; consequences are then determined by acts and states. The pertinent decision rule is simply maximizing expected utility. However, this separation is often not feasible, and the more general picture put forward by Fishburn (1964) is that everything is probabilistically assessed (except perhaps the acts themselves), though only conditionally on the possible acts. In this general picture, maximizing *conditional* expected utility is the appropriate decision rule. It may seem surprising that this situation in decision theory has so far not been reflected in equilibrium theory.

But this is not astonishing at all; that is my third point. The idea behind the general picture is that the conditional probabilities somehow hide causal dependencies which are more generally modeled in a probabilistic and not in a deterministic way (as Savage 1954 did). In the light of this idea, dependency equilibria are a mystery. The causal independence of

the choices of the players seems to be a defining characteristic of games in normal form. If Bob chooses after observing what Ann has chosen, then, of course, we have a clear case of a one way causal dependence. But how can Ann's choice then depend on Bob's? That would amount to a causal loop, and dependency equilibria seem to assume just this impossibility. However, the case is not as hopeless as it seems, as I shall try to indicate in Section 5. For the time being, let us look a little more closely at the properties of dependency equilibria.

3. Some Examples. The computation of dependency equilibria seems to be a messy business. Obviously it requires one to solve quadratic equations in two-person games, and the more persons, the higher the order of the polynomials we become entangled with. All linear ease is lost. Therefore, I cannot offer a well developed theory of dependency equilibria. Let us instead look at some much discussed simple games in order to develop a feeling for the new equilibria, namely, Matching Pennies, Bach or Stravinsky (BoS), Hawk and Dove, and PD. This discussion becomes more vivid when we consider the other kinds of equilibria for comparison. Afterwards, we can infer some simple theorems from these examples.

Matching Pennies. This is the paradigm for a pure conflict, that is, a zero-sum or constant-sum game. It is characterized by the following utility matrix:

| | | | |
|-------|-----|-------|-------|
| u | v | b_1 | b_2 |
| a_1 | | 0 | 1 |
| a_2 | | 1 | 0 |
| | | 0 | 1 |

It is clear that it has exactly one Nash equilibrium and exactly one correlated equilibrium. It is characterized by the following distribution:

| | | |
|-------|-------|-------|
| p | b_1 | b_2 |
| a_1 | 1/4 | 1/4 |
| a_2 | 1/4 | 1/4 |

By contrast, it is easily verified that the dependency equilibria of this game may be biased toward the diagonal or toward the counter-diagonal:

| | | |
|-------|-----------|-----------|
| p | b_1 | b_2 |
| a_1 | x | $1/2 - x$ |
| a_2 | $1/2 - x$ | x |

where $0 \leq x \leq 1/2$. It is instructive to represent the players' expected utilities in the various equilibria by a joint diagram (Figure 1).

Bach or Stravinsky. This game is a paradigmatic coordination game superimposed by a conflict. Its utility matrix is:

| | | |
|-------|-------|-------|
| v | b_1 | b_2 |
| a_1 | 1 | 0 |
| a_2 | 0 | 2 |

As is well known, this game has three Nash equilibria, two in pure strategies (the players can meet on the diagonal) and a mixed one:

| | | |
|-------|-------|-------|
| p | b_1 | b_2 |
| a_1 | 1 | 0 |
| a_2 | 0 | 0 |

| | | |
|-------|-------|-------|
| p | b_1 | b_2 |
| a_1 | 0 | 0 |
| a_2 | 0 | 1 |

| | | |
|-------|-------|-------|
| p | b_1 | b_2 |
| a_1 | 2/9 | 4/9 |
| a_2 | 1/9 | 2/9 |

The correlated equilibria of this game form just the convex closure of the Nash equilibria:

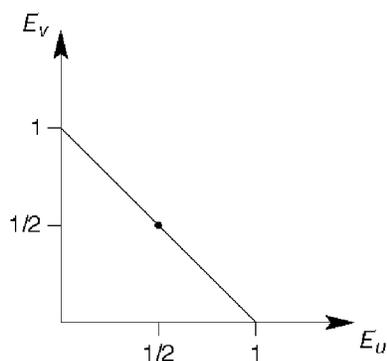


Figure 1. ‘·’ indicates Nash and correlated equilibria, and ‘-’ indicates dependency equilibrium.

| | | |
|-------|-----------------|---------------|
| p | b_1 | b_2 |
| a_1 | x | $\leq 2x, 2y$ |
| a_2 | $\leq x/2, y/2$ | y |

The dependency equilibria are again of three kinds:

| | | | | | |
|-------|-------|-------|-------|-------|-------|
| p | b_1 | b_2 | p | b_1 | b_2 |
| a_1 | 1 | 0 | a_1 | 0 | 0 |
| a_2 | 0 | 0 | a_2 | 0 | 1 |

provided the zero rows and columns are approximated in an appropriate way, and

| | | |
|-------|-----------|-----------|
| p | b_1 | b_2 |
| a_1 | x | $2/3 - x$ |
| a_2 | $1/3 - x$ | x |

where $0 \leq x \leq 1/3$. The players' expected utilities in these equilibria are shown in Figure 2. Quite similar observations can be made about pure coordination games without conflict, like meeting at one of two places.

Hawk and Dove. This game represents another very frequent type of social situation. It will show even more incongruity among the equilibrium concepts. So far, one may have thought that the correlated equilibria are the convex closures of the Nash equilibria. But this is not true. I shall consider the utility matrix preferred by Aumann because it illustrates that there are correlated equilibria that Pareto-dominate mixtures of Nash

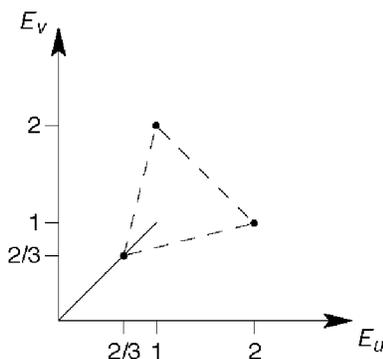


Figure 2. ‘•’ indicates Nash equilibrium, ‘-•-’ indicates correlated equilibrium, and ‘•-’ indicates dependency equilibrium.

equilibria; hence, both players may improve by turning to correlated equilibria. However, they may improve even more by looking at dependency equilibria. Here is the utility matrix:

| | | | |
|-------|-----|-------|-------|
| | v | | |
| u | | b_1 | b_2 |
| a_1 | | 6 | 7 |
| | | 6 | 2 |
| a_2 | | 2 | 0 |
| | | 7 | 0 |

There are again three Nash equilibria with the following expected utilities (see also Figure 3):

| | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| p | b_1 | b_2 | p | b_1 | b_2 | p | b_1 | b_2 |
| a_1 | 0 | 1 | a_1 | 0 | 0 | a_1 | 4/9 | 2/9 |
| a_2 | 0 | 0 | a_2 | 1 | 0 | a_2 | 2/9 | 1/9 |

The correlated equilibria reach out further on the diagonal (again see Figure 3). They are given by

| | | | |
|-------|--|-------|-------|
| p | | b_1 | b_2 |
| a_1 | | x | y |
| a_2 | | z | w |

where $x + y + z + w = 1$ and $0 \leq x/2, 2w \leq y, z$, and they yield the expected utilities shown in Figure 4.

Again, we have three kinds of dependency equilibria:

| | | | | | |
|-------|-------|-------|-------|-------|-------|
| p | b_1 | b_2 | p | b_1 | b_2 |
| a_1 | 0 | 1 | a_1 | 0 | 0 |
| a_2 | 0 | 0 | a_2 | 1 | 0 |

provided the zero rows and columns are approximated in an appropriate way, and

| | | | |
|-------|--|-------|--------------|
| p | | b_1 | b_2 |
| a_1 | | x | y |
| a_2 | | y | $1 - x - 2y$ |

where $y = (1/18)(2 - 15x + \sqrt{4 + 156x + 9x^2})$. This makes evident that we slip into quadratic equations. The corresponding expected utilities

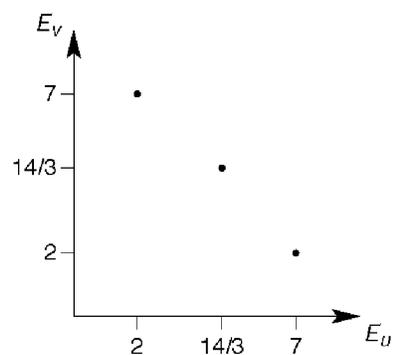


Figure 3.

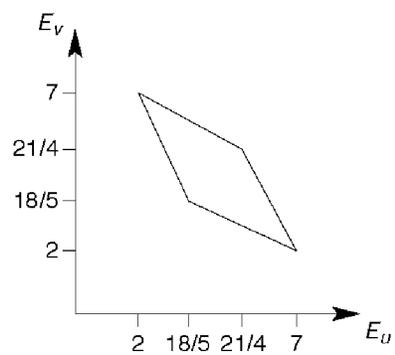


Figure 4.

reach out still further on the diagonal (see Figure 5). Clearly $6 > 21/4 > 14/3$, the maximal values reached on the diagonals of the three diagrams.

Prisoners' Dilemma. This is my final and perhaps most important example. Its utility matrix is

| u | v | b_1 | b_2 |
|-------|-----|-------|-------|
| a_1 | | 2 | 3 |
| a_2 | | 0 | 1 |
| | | 3 | 1 |

There is only one Nash equilibrium:

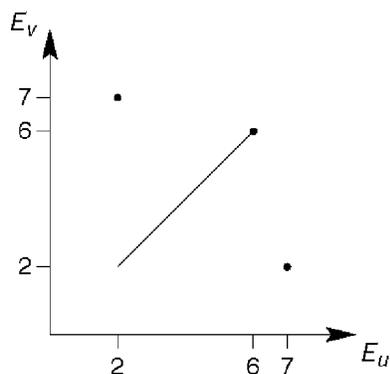


Figure 5.

| p | b_1 | b_2 |
|-------|-------|-------|
| a_1 | 0 | 0 |
| a_2 | 0 | 1 |

Indeed, defection ($= a_2$ or, respectively, b_2) strictly dominates cooperation ($= a_1$ or b_1); hence, there can be no other Nash equilibrium. For the same reason, this is also the only correlated equilibrium.

The dependency equilibria, by contrast, have a much richer structure. They come in two kinds:

| p | b_1 | b_2 |
|-------|---------------|-------------------|
| a_1 | $(1/2)x(1+x)$ | $(1/2)x(1-x)$ |
| a_2 | $(1/2)x(1-x)$ | $(1/2)(1-x)(2-x)$ |

where $0 \leq x \leq 1$, and

| p | b_1 | b_2 |
|-------|--------------------|--------------------|
| a_1 | $(3/8)(1-x)(1+x)$ | $(1/8)(1-x)(1-3x)$ |
| a_2 | $(1/8)(1+x)(1+3x)$ | $(3/8)(1-x)(1+x)$ |

where $-1/3 \leq x \leq 1/3$. The expected utilities in all these equilibria look very simple (see Figure 6).

It is of particular interest here that joint cooperation is among the dependency equilibria; indeed it weakly Pareto-dominates all other such equilibria. Of course, it is a well worn and very simple observation that such dependence between the players may make them cooperate. But now

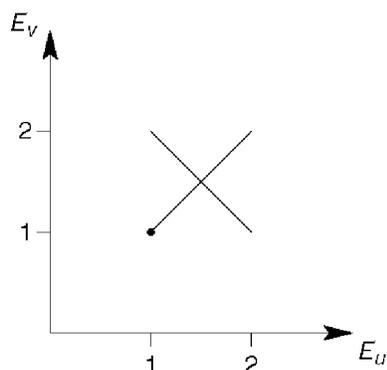


Figure 6. ‘•’ indicates Nash and correlated equilibria, and ‘-’ indicates dependency equilibrium.

we have found an equilibrium concept that underpins this observation. Moreover, we have seen that correlated equilibria do not provide the right kind of dependence for this purpose—they succumb to defection. Evidently, all this is strong motivation to try to make good sense of dependency equilibria.

4. Some Observations. The examples suggest some simple generalizations, all of which can be extended, it seems, to the n -person case.

Observation 1. Each Nash equilibrium of a two-person game is a correlated equilibrium. (*Proof:* Just look at the definitions.)

Observation 2. The set of correlated equilibria of a two-person game is convex.

Again, the proof is evident from the definition. Of course, we find both observations already in Aumann (1974, Section 4). They entail that the convex closure of the Nash equilibria of a game is a subset of the set of correlated equilibria.

The next observations are closer to our concerns:

Observation 3. Each Nash equilibrium of a two-person game is a dependency equilibrium. (*Proof:* Again, just look at the definitions.)

Observation 4. Generally, dependency equilibria are not included among the correlated equilibria, and vice versa. (*Proof:* Just look at the examples above.)

In BoS we saw that there are also very bad dependency equilibria, and in PD we luckily found one dependency equilibrium weakly Pareto-dom-

inating all the others. This suggests the following question: Which dependency equilibria are Pareto-optimal within the set of dependency equilibria? Clearly, these are the most interesting or attractive ones. Here is a partial answer:

Observation 5. Let $q = s \otimes t$ be a Nash equilibrium and suppose that the pure strategy combination (a_i, b_k) is at least as good as this equilibrium, that is, that $u_{ik} \geq \sum_{j,l} s_j t_l u_{jl}$ and $v_{ik} \geq \sum_{j,l} s_j t_l v_{jl}$. Then this combination, or p with $p_{ik} = 1$, is a dependency equilibrium.

Proof. Define $p^r = [(r-1)/r]p + (1/r)q$, and assume that p^r is strictly positive. Obviously $\lim_{r \rightarrow \infty} p^r = p$. Moreover, $\lim_{r \rightarrow \infty} \sum_l p_{li}^r u_{il} = u_{ik}$, and for all $j \neq i$ and all r , $\sum_l p_{lj}^r u_{jl} = \sum_l t_l u_{jl}$. But now we have $u_{ik} \geq \sum_{j,l} s_j t_l u_{jl} \geq \sum_l t_l u_{jl}$: the first inequality holds by assumption, and the second, because $\langle s, t \rangle$ is a Nash equilibrium. The same considerations apply to the other player. Hence, given our assumption, p with $p_{ik} = 1$ is a dependency equilibrium.

If p^r should not be strictly positive, modify q such that those a_j with $j \neq i$ and $s(a_j) = 0$ receive some positive probability by q , and such that $q(b_l|a_j) = t_l$, and correspondingly for those b_l with $l \neq k$ and $t(b_l) = 0$. Then the modified p^r is strictly positive, and the same proof goes through.

■

In PD, Hawk and Dove, and BoS this observation fully satisfies the quest for the Pareto-optima among the dependency equilibria. But it does not generally do so. In Matching Pennies no pure strategy combination is Pareto-better than the Nash equilibrium; yet mixtures of them in which equivalent strategy combinations have equal weight are dependency equilibria.

5. Discussion. As I have already indicated, dependency equilibria seem to be causal nonsense. Normal form games seem to be characterized by the causal independence of the actions of the players; in extensive form games, one action can influence the other, but there cannot exist causal loops, as apparently assumed by dependency equilibria. This objection is, however, already undermined by Reichenbach's common cause principle, according to which such a probabilistic correlation may always point to some common cause. So, in order to get clear about this point, we would have to engage in a most careful discussion of the causal structure of decision and game situations.

We cannot do this in the short space given here. Let me only indicate a few points. It is a most natural assumption that an action is caused by the decision situation of the agent. This is an internal state of the agent, her view of the situation consisting of her beliefs and desires, intentions, etc. Decision and game theory model such decision situations. The whole

model then represents the situation, but the situation itself is not an item in the model. It can be made so only in more complex models that I call 'reflexive decision models', something started in Eells (1982) and usefully applied to Newcomb's problem, but not further developed to my knowledge. Only in such reflexive models can the causation of the agent's actions be made explicit. They allow for a most useful separation of decision and action, they allow considering possible side effects of a decision situation besides the ensuing action, which may crucially matter to the decision, and they allow for a useful theory of commitment and a unification of sophisticated and resolute choice in the sense of McClennen (1990).⁵

Such considerations can be used, for example, to make a case for one-boxing in Newcomb's problem (from a causal point of view!). This already suggests a significance for PD, which is well known to be conceivable as a two-sided Newcomb problem. The point then is to conceive the decision situations of the players as somehow jointly caused and as entangled in a dependency equilibrium. The players are free to break the entanglement (in which case they defect) or to maintain the entanglement (in which case they are able to rationalize cooperation). However, by 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. This, in any case, would be my line of defense, if space would permit it.

So many pages have been filled with PD, and many at least resembling mine in spirit. So, let me add some comparative remarks, which may shed some further light on the new equilibrium concept.

1. Some philosophers may claim to have offered a much simpler rationalization of cooperation in the one-shot PD than the one I have put forward, namely, via the so called mirror principle (cf. Davis 1977 and Sorensen 1985), which says that whenever Ann and Bob are in the same decision situation, they act in the same way. In PD they are in the same situation because of the symmetry of the story. Hence, only joint cooperation and joint defection are possible outcomes. If they believe in the principle, they also believe that these are the only possible outcomes. Hence, since for both joint cooperation is better than joint defection, it is rational for both to cooperate.

This argument is entailed by my account, so I agree with its conclusion. But I find it too quick. It avoids causal considerations, and it does not present a theory of rationality that would entail the rationality of cooperation for Ann and Bob in their particular situation. Therefore, it does not exclude that mutual defection satisfies the mirror principle as well, as

5. All of this is more fully explained in Spohn 2003.

the standard theory has it. In a way, it takes its conclusion for granted. By contrast, my account attempts to back up the mirror argument by specifying a theory of rationality in which the rationality of cooperation emerges as a conclusion. Moreover, it is hard to see how to apply the mirror principle when the situation is not symmetric; but there is no such presupposition in dependency equilibria.

2. In a similar spirit, some game theorists may say that my account really belongs to *cooperative game theory*, within which the cooperative solution of PD is no mystery anyway. However, if we follow Osborne and Rubinstein (1994, 255ff.) and take cooperative game theory to refer to groups of players without considering “the details of how groups of players function internally,” this is not true; dependency equilibria are intended to rationalize cooperation as an account of individual rationality. If we follow Harsanyi and Selten (1988, 1) instead and define “cooperative games as those permitting enforceable agreements,” it is again not true; my rationalization of cooperation refrained from alluding to enforceable agreements. Whether it will succeed is, of course, another question. In any case, dependency equilibria may provide a story about individual rationality backing up cooperative game theory to some extent.

3. To continue on the issue, Harsanyi and Selten (1988, 4–7) showed how cooperation can emerge as an ordinary Nash equilibrium in a non-cooperative PD by adding a preplay of commitment moves. Ann starts making her conditional commitment move (“I commit myself to cooperate, if Bob does so as well”); Bob follows with his unconditional commitment move (“I commit myself to cooperate”); and then both cooperate. Thereby we do not even have to assume a causal *interdependence* of the players’ decision situations; the dependence is successively generated by the commitment moves. But, of course, the idea of Harsanyi and Selten is that there is some external mechanism sufficiently lowering the payoffs in case the commitments are violated. By contrast, the indicated account of commitment provided by reflexive decision theory hopes to do without such external sanctions.

4. My proposal closely resembles the old theory of conjectural variations about duopoly that is well reviewed in Friedman 1983, Sections 5.1–5.2). He concludes (107) that “at the level of simultaneous decisions in a single period model, conjectural variation is not meaningful,” and hence interprets the old single period models of conjectural variation as being implicitly about a dynamical process that has been more recently treated in multiperiod models. So, does this verdict apply as well to the account offered here? Yes, at least according to received view: the effective objection has been that single period conjectural variations assume a causal dependence that simply does not exist. However, the longer story I have indicated above tries to dispel exactly this objection. Whether my argu-

ment would help making sense of the theory of conjectural variations is, however, beyond my judgment.

5. I would like at least to mention the work of Albert and Heiner (2003) and Heiner, Albert, and Schmidtchen (2000).⁶ They also seek a causally unassailable rationalization of one boxing in Newcomb's problem (something I have only mentioned here) and proceed to generalize their account to a treatment of the one-shot PD. There are obvious differences, due to the fact that, *prima facie*, the setting of their story is evolutionary game theory, which provides quite a different frame of interpretation. Still, the similarities in intention and procedure are quite remarkable.

REFERENCES

- Albert, M., and R. A. Heiner (2003), "An Indirect-Evolution Approach to Newcomb's Problem", *Homo Oeconomicus* 20: 161–194.
- Aumann, R. (1974), "Subjectivity and Correlation in Randomized Strategies", *Journal of Mathematical Economics* 1: 67–96.
- (1987), "Correlated Equilibrium as an Expression of Bayesian Rationality", *Econometrica* 55: 1–18.
- (1995), "Backward Induction and Common Knowledge of Rationality", *Games and Economic Behavior* 8: 6–19.
- Davis, L. (1977), "Prisoners, Paradox, and Rationality", *American Philosophical Quarterly* 114: 319–327.
- Eells, E. (1982), *Rational Decision and Causality*. Cambridge: Cambridge University Press.
- Fishburn, P. C. (1964), *Decision and Value Theory*. New York: Wiley.
- Friedman, J. W. (1983), *Oligopoly Theory*. Cambridge: Cambridge University Press.
- Harsanyi, J. C., and R. Selten (1988), *A General Theory of Equilibrium Selection in Games*. Cambridge, MA: MIT Press.
- Heiner, R. A., M. Albert, and D. Schmidtchen (2000), "Rational Contingent Cooperation in the One-Shot Prisoner's Dilemma", manuscript.
- McClellenn, E. F. (1990), *Rationality and Dynamic Choice*. Cambridge: Cambridge University Press.
- Myerson, R. B. (1991), *Game Theory: Analysis of Conflict*. Cambridge, MA: Harvard University Press.
- Osborne, M. J., and A. Rubinstein (1994), *A Course in Game Theory*. Cambridge, MA: MIT Press.
- Rabinowicz, W. (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*, Lund University Reports 2000:1. Lund: Lund University, 243–256.
- Savage, L. J. (1954), *The Foundations of Statistics*. New York: Dover.
- Sorensen, R. A. (1985), "The Iterated Versions of Newcomb's Problem and the Prisoner's Dilemma", *Synthese* 63: 157–166.
- Spohn, W. (1982), "How to Make Sense of Game Theory", in W. Stegmüller, W. Balzer, and W. Spohn (eds.), *Philosophy of Economics*. Berlin: Springer, 239–270.
- (2003), "Dependency Equilibria and the Causal Structure of Decision and Game Situations", *Homo Oeconomicus* 20: 195–255.

6. See also the other papers by Heiner in the CSLE Discussion Paper Series of the Center for the Study of Law and Economics at the Universität des Saarlandes.