

Cognitive Uncertainty in Games

A Note on Limited Information Processing and Backward Induction

Thorsten Clausing *

*Faculty of Economics and Management, University of Magdeburg, Postfach 4120,
39016 Magdeburg, Germany*

E-mail: Thorsten.Clausung@ww.uni-magdeburg.de

A notion of cognitive uncertainty is introduced as an agent's uncertainty about the validity of the results of his own information processing. In order to analyze this notion, a formal model of the agent's information processing is presented. It is shown how cognitive uncertainty may prevent a group of well informed rational agents from playing backward induction in a game of perfect information, whereas they would play backward induction without such uncertainty.

Key Words: uncertainty, information processing, backward induction

1. INTRODUCTION

In economic theory, the modeling of uncertainty plays an important role. Most notably this topic is treated in the theory of expected utility maximization and in Bayesian game theory. In these contexts, uncertainty is usually taken to arise as the result of incomplete or unreliable information about the decision problem faced by an agent. The agent uses this information to form expectations about the relative likelihood of different possible consequences of the actions from which he has to choose and bases his decision on these estimates. However, when one takes such an information processing interpretation of decision making seriously, one may identify a second possible source of uncertainty. Even when an agent perceives his initial information as complete and reliable, he may not be convinced of the correctness of the results of his information processing because he may fear that he made reasoning mistakes in drawing conclusions from the ini-

* I benefited from comments by Ulrich Berger, Dominique Demougin, Atanasios Mitropoulos, Arnis Vilks, and seminar participants at the University of Magdeburg, the University of Leipzig and the Bonn Summer School on Interactive Economic Decisions.

tial information. As an illustration, consider the situation of a chess player who has just found a combination that forces a checkmate in nine moves by sacrificing his queen. Now this player certainly knows the position of the pieces on the board as well as the rules of the game, so there can be no uncertainty about the structure of his decision problem. Nevertheless, he may very well hesitate to sacrifice his queen because he has doubts as to the correctness of his combinatorial calculations. I will refer to uncertainty of this kind as cognitive uncertainty, in contrast to the uncertainty stemming from defective information which I will call informational uncertainty.

The aim of this paper is to study the consequences of cognitive uncertainty in games. Attention is restricted to games of perfect information that can be unambiguously solved by backward induction in the sense that at all decision nodes, there is a unique payoff maximizing move under the assumption that the respective backward induction moves will be played at all later nodes. In order to analyze cognitive uncertainty formally, I will develop a model that describes the agent's information processing explicitly. It is supposed to capture a situation where the agent has complete and reliable information about the structure of the game that is going to be played. He also knows that his opponents are rational and have the same information as he himself.

A word on the notion of rationality is in order here. Most economists would probably call an agent who is not convinced of the correctness of his own information processing boundedly rational. Hence one might call this paper a note on bounded rationality and backwards induction. However, I prefer to use the term rationality as in the literature on epistemic foundations of game theory (see, e.g., Dekel and Gul [1]), where it just refers to the relation between what the agents do and what they know or believe. Thus in my terminology, an agent can be rational even though his information processing is not perfect.

The remainder of this paper is organized as follows. In section 2, an intuitive informal as well as a formal presentation of the information processing model are given. Section 3 contains the results on the choice of moves with cognitive uncertainty. Section 4 contains a further discussion of the theory of cognitive uncertainty and compares it to the more usual approach employing informational uncertainty. Proofs are given in section 5.

2. THE MODEL

Let me start with an informal and somehow simplified description of the information processing model. Its basic idea is to describe how the agent's set of beliefs, i.e. his mental "picture of the world" evolves while he processes his initial information before the start of the game. To this end,

the information processing is broken down into a sequence of successive steps. By looking at the belief sets at the completion of each of these steps, one thus has a sequence of "mental pictures" at consecutive points of time. The first picture just reflects the initial information, and each succeeding one in the sequence results from applying one more step of information processing to the preceding picture.

I will represent these pictures with the help of sets of states. Intuitively, each state stands for a future situation that the agent considers possible, i.e. for one way how the decision problem he faces might evolve. Of course, in the given context a future situation is just a play of the given game that arises, and one can therefore identify each state with one of the plays of this game. That an agent considers a situation possible then just means that at the present stage of his information processing, he is not convinced that this play will not actually arise.

In order to capture the idea that the agent is completely informed about the structure of the given game, the set of states representing the initial information consists of exactly one state for each of its possible plays. The less states are contained in the representation of a mental picture, the sharper the picture; therefore the set of states should shrink during the information processing.

The idea that the agent knows that all players behave rationally is captured by having him at each information processing step delete all those plays from his set of states that a rational player would not choose. Rational behavior here means that a player will not decide to take a certain move if he knows that an alternative move gives him a higher expected payoff. Of course, what a player knows about the expected payoffs of different moves depends upon how many steps of information processing he has already completed. A given move may well seem acceptable after three steps but lose this property after the fifth step. I will therefore introduce a notion of t -irrationality as the property of a move that is perceived as inferior to an alternative move and thus not taken by a player who has completed at least t steps of information processing. Note that by iteratively deleting moves of his opponents, the agent decides on the basis of what he has derived after t steps of *his* information processing which moves his rational opponents will decide not to play after t steps of *their* information processing. I thus implicitly assume that the agent not only knows that his opponents have the same initial information as he himself, but also assumes that they process this information in the same way as he does.

Finally, cognitive uncertainty now enters the stage when the expected payoff of a move is calculated. This obviously calls for a probability measure on the set of states. If the agent were certain that he processed his information correctly, i.e. that the plays which he deleted from his set of states will indeed not arise, then he should leave these plays out of con-

sideration when determining the expected payoff of a move. Consequently, a probability measure should be used that assigns probability zero to any deleted state. However, with cognitive uncertainty, these states cannot be completely ignored, and they could consequently be assigned some positive probability. I introduce an upper bound on the probability that may be assigned to deleted states. This bound increases in the number of processing steps to capture the idea that information processing results are the less reliable the more processing steps they demand. The maximum probability that may be assigned to deleted states as the number of steps goes to infinity yields an uncertainty parameter $\epsilon \in [0, 1]$ which can be taken as a natural measure of the degree of cognitive uncertainty of a given agent.

It is now time to give a formal account of the model. To this end, I introduce the following notation. Let the perfect information game Γ be played. A possible move X at decision node v of Γ is represented as X_v and the set of its possible moves as \mathcal{M} . Σ denotes the set of possible plays of Γ , $m = \#\Sigma$ the number of possible plays and $[X_v] \subset \Sigma$ the set of all plays containing the move X_v . For a play $\sigma \in \Sigma$, $\pi_i(\sigma)$ stands for the payoff received by player i if σ obtains. Furthermore, $i(v)$ denotes the player having to move at decision node v and $\mathcal{P}(X)$ the power set of X .

The model consists of a sequence $(\Sigma_t, \mathcal{P}_t, \mathcal{D}_t)_{t \in N}$. For all natural numbers t , Σ_t is a subset of Σ and \mathcal{P}_t is a set of probability measures on $(\Sigma, \mathcal{P}(\Sigma))$. The set \mathcal{D}_t contains vectors of distance functions defined as follows. For any play $\sigma \in \Sigma$, call a bijection $d_\sigma : \Sigma \rightarrow \{0, \dots, m-1\}$ a distance function if it satisfies condition (1):

$$d_\sigma(\sigma) = 0 \tag{1}$$

A vector of distance functions $d = (d_{\sigma_1}, \dots, d_{\sigma_m})$ contains exactly one distance function for every play $\sigma \in \Sigma$. These functions are needed to introduce a counterfactual element into the model. For an intuitive interpretation, take $d_\sigma(\sigma')$ to denote the distance or dissimilarity between σ and σ' . With this interpretation, condition (1) appears quite natural. Clearly, a player i reasoning about what moves to choose should consider counterfactual sentences of the kind "If move X_v were made, I would receive a payoff of x ", and thus the use of a counterfactual element in the model appears well motivated. Following Lewis [3], I interpret such a counterfactual sentence to mean that in a hypothetical world in which X_v is played, but which otherwise is as similar to the actual world as possible, player i receives a payoff of x . For a world where play σ actually arises, the play σ' corresponding to this hypothetical world can now be determined with the help of the vector d : it is just the state in $[X_v]$ at which d_σ takes

its minimum value. Thus one can define a counterfactual payoff function $\pi_d : \Sigma \times \mathcal{M} \rightarrow \mathbb{R}$ for any move $X_v \in \mathcal{M}$ as follows:

$$\pi_d(\sigma, X_v) = \pi_{i(v)}(\sigma') \text{ such that } \sigma' = \underset{\hat{\sigma} \in [X_v]}{\operatorname{argmin}} d_\sigma(\hat{\sigma}) \quad (2)$$

As d_σ is injective on a finite domain, argmin is a singleton set which I identify with its only element in (2) for notational simplicity. Now for each play σ , the function π_d associates with each move X_v a payoff for the player who moves at v . Note that if X_v belongs to σ , this is just the payoff received if σ obtains because of (1). Otherwise, as described above, it is the payoff that is received if the play with the least distance to σ obtains. Of course, in general such counterfactual payoffs are not uniquely determined by the game tree of Γ . However, for any counterfactual payoff that player $i(v)$ could possibly receive from choosing X_v , there is a vector of distance functions d generating it.

The sequence $(\Sigma_t, \mathcal{P}_t, \mathcal{D}_t)_{t \in \mathbb{N}}$ can now be defined inductively as follows. For $t = 1$, one has $\Sigma_1 = \Sigma$. \mathcal{D}_1 is the set of all distance functions and \mathcal{P}_1 is the set of all probability measures on $(\Sigma, \mathcal{P}(\Sigma))$.

Now assume $(\Sigma_t, \mathcal{P}_t, \mathcal{D}_t)$ has been defined. For a given probability measure $p \in \mathcal{P}_t$ and a given vector of distance functions $d \in \mathcal{D}_t$ it is now possible to calculate the expected payoff of a move X_v as follows.

$$E_{p,d}(\pi_{i(v)}|X_v) = \sum_{\sigma \in \Sigma} p(\sigma) \cdot \pi_d(\sigma, X_v) \quad (3)$$

I call a move X_v t -rational if it belongs to some play $\sigma \in \Sigma_t$ and it is not the case that there is a move Y_v belonging to some play $\sigma' \in \Sigma_t$ such that for all $d \in \mathcal{D}_t$ and all $p \in \mathcal{P}_t$, the expected payoff of Y_v is higher than that of X_v . Any move X_v that is not t -rational will be called t -irrational. A play is t -rational if it consists only of t -rational moves.

Now $(\Sigma_{t+1}, \mathcal{P}_{t+1}, \mathcal{D}_{t+1})$ can be defined as follows. Let Σ_{t+1} contain exactly the t -rational plays $\sigma \in \Sigma_t$. Furthermore, let the set \mathcal{D}_{t+1} consist of all distance functions that satisfy the following condition.

$$d_\sigma(\sigma') < d_\sigma(\sigma'') \text{ if } \sigma, \sigma' \in \Sigma_{t+1}, \sigma'' \in \Sigma \setminus \Sigma_{t+1} \quad (4)$$

Intuitively, restriction (4) says that any undeleted play is less distant from any other undeleted play than from any deleted play.

Finally define the set of allowable probability measures \mathcal{P}_{t+1} to contain all probability measures p on $(\Sigma, \mathcal{P}(\Sigma))$ that satisfy the following restriction for a given constant $\epsilon \in [0, 1]$.

$$p(\Sigma_{t+1}) \geq 1 - \frac{t}{t+1}\epsilon \quad (5)$$

As described above, without cognitive uncertainty, the agent should assign zero probability to any deleted play $\sigma \in \Sigma \setminus \Sigma_{t+1}$. The weaker condition (5), on the contrary, allows him to assign some positive probability to such plays, reflecting his doubts about the results of his information processing. For the reasons stated above, condition (5) is chosen such that the maximum probability assignable to deleted plays increases in t and at the limit for $t \rightarrow \infty$, this maximum probability goes up to ϵ .

A remark on the sets \mathcal{P}_t and \mathcal{D}_t is in order here. To work with sets of probability measures instead of a single measure for each agent may seem somewhat unusual. However, these sets intuitively contain all probability assignments that are compatible with the agent's state of knowledge before his t -th step of information processing has been completed. Taking into consideration all of these compatible probability measures when thinking about which moves his opponents will not play can be seen as an act of caution by the agent. An analogous statement can be made about the sets of vectors of distance functions.

Let me conclude this section with the following lemma which says that even though Σ_t weakly shrinks as t becomes larger, the agent always considers at least one play of the game as possible. This may be seen as a consistency requirement on his picture of the world.

LEMMA 2.1. *For all natural numbers t , $\Sigma_t \neq \emptyset$.*

3. RESULTS

As a reference case, consider a situation without cognitive uncertainty, i.e. $\epsilon = 0$. Let m denote the maximum number of moves in any play of Γ . For this case it can be shown that the player will come to believe that the backward induction play, which I denote by σ_{BI} , will arise.

THEOREM 3.1.

For $\epsilon = 0$ and $t > m$, $\Sigma_t = \{\sigma_{BI}\}$

The proof of this theorem makes use of the following lemma, which I state here for latter reference. Define $[X_v]_t = [X_v] \cap \Sigma_t$. For any move X_v such that $[X_v]_t \neq \emptyset$, define $\pi_t^{Min}(X_v)$ to denote the lowest payoff the player $i(v)$ gets in any play $\sigma \in \Sigma_t$ where he plays X_v , i.e. $\pi_t^{Min}(X_v) = \min\{\pi_{i(v)}(\sigma) \mid \sigma \in [X_v]_t\}$. $\pi_t^{Max}(X_v)$ is defined accordingly with max instead of min.

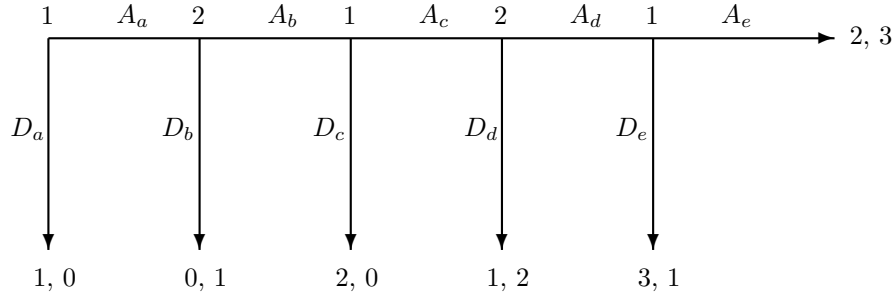


FIG. 1. A game of perfect information

Furthermore, let $\pi^{Min}(X_v)$ stand for the lowest payoff player $i(v)$ gets in any play $\sigma \in \Sigma$ where he plays X_v , i.e. $\pi^{Min}(X_v) = \min\{\pi_{i(v)}(\sigma) \mid \sigma \in [X_v]\}$, and $\pi^{Max}(X_v)$ is again defined accordingly with max instead of min.

LEMMA 3.1. *Let X_v be a move belonging to some state $\sigma \in \Sigma_t$. Then X_v is t -irrational if and only if there is some move Y_v belonging to some play $\sigma' \in \Sigma_t$ such that:*

$$(1 - \frac{t-1}{t}\epsilon)\pi_t^{Min}(Y_v) + \frac{t-1}{t}\epsilon\pi^{Min}(Y_v) > (1 - \frac{t-1}{t}\epsilon)\pi_t^{Max}(X_v) + \frac{t-1}{t}\epsilon\pi^{Max}(X_v)$$

With $\epsilon > 0$, however, the backward induction result from theorem 3.1 no longer holds. To see this, consider the centipede-style game in Fig. 1. Its backward induction play consists of the single move D_a .

As there are six possible plays of this game, Σ has six elements. I will refer to the play containing A_e as σ_6 , the play containing D_e as σ_5 , and all other plays are indexed by the number of moves belonging to them. By definition, $\Sigma_1 = \Sigma$. Because of $\pi^{Min}(D_e) = \pi_1^{Min}(D_e) = 3$ and $\pi^{Max}(A_e) = \pi_1^{Max}(A_e) = 2$, lemma 3.1 shows that A_e is 1-irrational and thus σ_6 does not belong to Σ_2 . One can easily check that all other moves are 1-rational, which means $\Sigma_2 = \{\sigma_i \mid 1 \leq i \leq 5\}$.

Now check whether A_d is 2-rational. For this to be the case, lemma 3.1 says that the following inequality must not be fulfilled:

$$(1 - \frac{\epsilon}{2})\pi_2^{Min}(D_d) + \frac{\epsilon}{2}\pi^{Min}(D_d) > (1 - \frac{\epsilon}{2})\pi_2^{Max}(A_d) + \frac{\epsilon}{2}\pi^{Max}(A_d) \quad (6)$$

This inequality is fulfilled exactly if $\epsilon < 1$, which means that only in the extreme case $\epsilon = 1$ does σ_5 belong to Σ_3 . Note that if A_d is 2-rational, it is also t -rational for any $t > 2$. This is because in the inequality (6), the left hand side is independent of t because of $\pi_t^{Min}(D_d) = \pi^{Min}(D_d)$ for all natural numbers t and the right hand side increases in t . Again, one may easily check that all moves other than A_d are necessarily 2-rational. Now one can proceed in the same way to find that σ_4 does not belong to Σ_4 if and only if $\epsilon < 3/4$, σ_3 does not belong to Σ_5 if and only if $\epsilon < 4/9$, and σ_2 does not belong to Σ_6 if and only if $\epsilon < 5/12$. Finally, σ_1 belongs to Σ_t for any natural number t and any value of ϵ .

Thus sufficiently high cognitive uncertainty can cause a player with complete and reliable information about the game and the rationality of his opponents not to expect exclusively the backward induction play.

THEOREM 3.2. *With $\epsilon > 0$, there can be a play $\sigma \neq \sigma_{BI}$ such that $\sigma \in \Sigma_t$ for all natural numbers t .*

Note however that it can easily be seen from lemma 3.1 that for any node at which only terminal moves are possible, only the backward induction move can belong to any play $\sigma \in \Sigma_t$, $t \geq 2$.

4. DISCUSSION

The main interest in theorem 3.2 comes probably from the fact that in experiments with centipede-style games, the backward induction outcome is rarely observed (cf. McKelvey and Palfrey [4]). Taking cognitive uncertainty into account provides an explanation of these experimental results. However, deviations from backward induction may also be explained with the help of informational uncertainty, which is indeed the route usually taken in the literature (cf. Kreps et alii [2]; Zauner [5]). In these articles it is assumed that the players are from the outset uncertain about the structure of the game and/or the rationality of their opponents, and it is shown how with such beliefs it may be rational not to play backward induction. What then do we need a theory of cognitive uncertainty for?

Obviously there are many situations where economic agents have from the outset only incomplete information about the decision problem they face. Thus, for example, an auctioneer might know the average valuation of an item in the population of all potential buyers quite well, but still be uncertain about the particular group of buyers sitting in front of him at a given auction. However, there are also situations where it would seem that all agents involved in a given decision problem are endowed with all the relevant information. In particular in an experimental setting where all players are jointly instructed about the game to be played, there is little

reason to suspect that they might be uncertain about its structure. Also there seems to be little motivation for the players to expect that their opponents are not rational in the rather weak sense of not knowingly choosing dominated moves, as captured by the definition of t -rationality. Of course, even in such a setting there may be uncertainty as to how the monetary payoffs in the game translate into utility values. However, the kind of informational uncertainty assumed in the literature to rationalize deviations from backward induction is also often quite specific, demanding for example that players believe that there is a given percentage of opponents who will always play across in a centipede game. In short, there are situations where there seems to be only minimal informational uncertainty and the cognitive uncertainty explanation appears more natural.

It is, of course, in principle possible to reinterpret the uncertainty that is used in the rationalization literature on backward induction experiments as cognitive uncertainty, even though this is not the story usually told. The formalisms employed in these articles are independent of that story and one may just claim that the uncertainty somehow arose during the information processing which is treated as a black box. However, the advantage of the model presented here against such a reinterpretation is clearly that it provides an explicit description of the processing and thus an account of how cognitive uncertainty comes about.

Finally it is worth noting that by embedding it in a more complete model, one can derive some additional testable predictions from my information processing model which are not obtained with the rationalization approach. Thus imagine that there is a population of players with different individual ϵ -values. Let these ϵ -values be non-degenerately distributed on the interval $[0, 1]$. As seen in the example of the preceding section, the lower a player's ϵ -value, the more can be said about which plays will not arise. If you randomly select players from the population, this translates into a probability statement saying that the higher the supinum of ϵ -values for which a given play can be excluded, the less likely this play is to be observed. On this basis, one can construct e.g. two centipede-style games of equal length with different supinum ϵ -values for the exclusion of plays reaching the last nodes of the game tree. The conjecture that these nodes will be reached less frequently in the game with the higher supinum can then be checked experimentally.

Furthermore, even though ϵ has been introduced as an individually given constant above, one may extend the model and think of different factors that influence a player's ϵ and therefore the probability of observing backward induction play. For example, if you take ϵ to denote the thoroughness of information processing, one may conjecture that in games with high stakes players will pay more attention to how they ought to choose, i.e. process their information more thoroughly and are thus more likely to play

backward induction moves. Thus one could again construct two games, this time with the same game tree, but where all payoffs in the second game are those in the first game multiplied by some great constant, and check the conjecture that in the high payoff game backward induction behavior is more frequent experimentally.

5. PROOFS

Proof (of lemma 3.1). Let X_v and Y_v be such that $[X_v]_t \neq \emptyset$ and $[Y_v]_t \neq \emptyset$. Assume that the inequality in the statement of the lemma is not met. I will show that X_v is t -rational by constructing a vector of distance functions $d \in \mathcal{D}_t$ and a probability measure $p \in \mathcal{P}_t$ such that $E_{p,d}(\pi_{i(v)}|X_v) \geq E_{p,d}(\pi_{i(v)}|Y_v)$.

In order to construct a vector of distance functions $d \in \mathcal{D}_t$ and a probability measure $p \in \mathcal{P}_t$ with the necessary properties I distinguish the three mutually exclusive cases where $\pi^{Min}(Y_v) = \pi_t^{Min}(Y_v)$ and $\pi^{Max}(X_v) = \pi_t^{Max}(X_v)$, where $\pi^{Min}(Y_v) < \pi_t^{Min}(Y_v)$, and where $\pi^{Max}(X_v) > \pi_t^{Max}(X_v)$ and $\pi^{Min}(Y_v) = \pi_t^{Min}(Y_v)$.

Case 1: $\pi^{Min}(Y_v) = \pi_t^{Min}(Y_v)$ and $\pi^{Max}(X_v) = \pi_t^{Max}(X_v)$. By the definition of $\pi_t^{Min}(Y_v)$ and $\pi_t^{Max}(X_v)$, there are plays $\sigma \in [Y_v]_t$ and $\sigma' \in [X_v]_t$ such that:

$$\pi_{i(v)}(\sigma) = \pi_t^{Min}(Y_v) \quad \text{and} \quad \pi_{i(v)}(\sigma') = \pi_t^{Max}(X_v) \quad (7)$$

Define d_σ so that $d_\sigma(\sigma') = 1$. Then for any $d \in \mathcal{D}_t$ containing d_σ , one has

$$\pi_d(\sigma, Y_v) = \pi_t^{Min}(Y_v) \quad \text{and} \quad \pi_d(\sigma, X_v) = \pi_t^{Max}(X_v). \quad (8)$$

For the probability measure $p \in \mathcal{P}_t$, set $p(\sigma) = 1$ and $p(\sigma^*) = 0$ for all $\sigma^* \in \Sigma \setminus \{\sigma\}$.

Case 2: $\pi^{Min}(Y_v) < \pi_t^{Min}(Y_v)$. Note that this implies $(\Sigma \setminus \Sigma_t) \cap [Y_v] \neq \emptyset$. There are again plays $\sigma \in [Y_v]_t$ and $\sigma' \in [X_v]_t$ satisfying (7). Furthermore, there are plays $\sigma'' \in (\Sigma \setminus \Sigma_t) \cap [Y_v]$ and $\sigma''' \in [X_v]$ such that:

$$\pi_{i(v)}(\sigma'') = \pi^{Min}(Y_v) \quad \text{and} \quad \pi_{i(v)}(\sigma''') = \pi^{Max}(X_v) \quad (9)$$

Define d_σ as in case 1 and $d_{\sigma''}$ so that $d_{\sigma''}(\sigma''') = 1$. Then for any $d \in \mathcal{D}_t$ containing d_σ and $d_{\sigma''}$, one has (8) and

$$\pi_d(\sigma'', Y_v) = \pi^{Min}(Y_v) \quad \text{and} \quad \pi_d(\sigma'', X_v) = \pi^{Max}(X_v) \quad (10)$$

To construct $p \in \mathcal{P}_t$, let $p(\sigma) = 1 - \frac{t-1}{t}\epsilon$, $p(\sigma'') = \frac{t-1}{t}\epsilon$ and $p(\sigma^*) = 0$ for all $\sigma^* \in \Sigma \setminus \{\sigma, \sigma''\}$.

Case 3: $\pi^{Max}(X_v) > \pi_t^{Max}(X_v)$ and $\pi^{Min}(Y_v) = \pi_t^{Min}(Y_v)$. Here the same construction as in case 2 applies where Y_v and X_v , Min and Max , and $<$ and $>$ have to be exchanged one for another.

In all three cases, one finds for the expected payoff of X_v and Y_v with $d \in \mathcal{D}_t$ and $p \in \mathcal{P}_t$ as constructed:

$$\begin{aligned} E_{p,d}(\pi_{i(v)}|X_v) &= (1 - \frac{t-1}{t}\epsilon)\pi_t^{Max}(X_v) + \frac{t-1}{t}\epsilon\pi^{Max}(X_v) \geq \\ &(1 - \frac{t-1}{t}\epsilon)\pi_t^{Min}(Y_v) + \frac{t-1}{t}\epsilon\pi^{Min}(Y_v) = E_{p,d}(\pi_{i(v)}|Y_v) \end{aligned}$$

Now assume that the inequality in the statement of the lemma is met. For all $\sigma \in \Sigma$ and a given move Y_v with $[Y_v]_t \neq \emptyset$ one trivially has $\pi_d(\sigma, Y_v) \geq \pi^{Min}(Y_v)$ for any $d \in \mathcal{D}_t$. If furthermore $\sigma \in \Sigma_t$, by (4) the play with minimal distance to σ in $[Y_v]$ is some $\sigma' \in [Y_v]_t$. This means $\pi_d(\sigma, Y_v) \geq \pi_t^{Min}(Y_v)$ for any $d \in \mathcal{D}_t$.

Analogously, for a given move X_v with $[X_v] \neq \emptyset$ one has $\pi_d(\sigma, X_v) \leq \pi^{Max}(X_v)$ for all $\sigma \in \Sigma$ and $\pi_d(\sigma, X_v) \leq \pi_t^{Max}(X_v)$ for all $\sigma \in \Sigma_t$ and any $d \in \mathcal{D}_t$.

With any $p \in \mathcal{P}_t$ one thus finds:

$$\begin{aligned} E_{p,d}(\pi_{i(v)}|Y_v) &= \sum_{\sigma \in \Sigma_t} p(\sigma) \cdot \pi_d(\sigma, Y_v) + \sum_{\sigma \in \Sigma \setminus \Sigma_t} p(\sigma) \cdot \pi_d(\sigma, Y_v) \geq \\ &\sum_{\sigma \in \Sigma_t} p(\sigma) \cdot \pi_t^{Min}(Y_v) + \sum_{\sigma \in \Sigma \setminus \Sigma_t} p(\sigma) \cdot \pi^{Min}(Y_v) \geq \\ &(1 - \frac{t-1}{t}\epsilon)\pi_t^{Min}(Y_v) + \frac{t-1}{t}\epsilon\pi^{Min}(Y_v) \end{aligned} \quad (11)$$

Analogously one can derive:

$$E_{p,d}(\pi_{i(v)}|X_v) \leq (1 - \frac{t-1}{t}\epsilon)\pi_t^{Max}(X_v) + \frac{t-1}{t}\epsilon\pi^{Max}(X_v) \quad (12)$$

This gives $E_{p,d}(\pi_{i(v)}|X_v) < E_{p,d}(\pi_{i(v)}|Y_v)$ for all $p \in \mathcal{P}_t$ and $d \in \mathcal{D}_t$, which means that X_v is t -irrational. ■

Proof (of lemma 2.1). Assume that Σ_t is not empty, which is trivial for $t = 1$. I will show that Σ_{t+1} is not empty either. To this end, denote the set of moves that are played at node v in some play in Σ_t by $\mathcal{M}_t^v = \{X_v|[X_v]_t \neq \emptyset\}$. Define a binary relation \prec_t on \mathcal{M}_t^v such that $X_v \prec_t Y_v$

if and only if the inequality in the statement of lemma 3.1 is met. \prec_t is transitive: From $X_v \prec_t Y_v$ and $Z_v \prec_t X_v$ follows $Z_v \prec_t Y_v$ because of

$$\left(1 - \frac{t-1}{t}\epsilon\right)\pi_t^{Max}(X_v) + \frac{t-1}{t}\epsilon\pi_t^{Max}(X_v) \geq \left(1 - \frac{t-1}{t}\epsilon\right)\pi_t^{Min}(X_v) + \frac{t-1}{t}\epsilon\pi_t^{Min}(X_v) \quad (13)$$

Furthermore, (13) also implies that \prec_t is irreflexive. \mathcal{M}_t^v is clearly finite. If node v is reached in some play $\sigma \in \Sigma_t$, i.e. $\mathcal{M}_t^v \neq \emptyset$, one can thus show that a contradiction would arise if for all $X_v \in \mathcal{M}_t^v$ there were a move Y_v with $X_v \prec_t Y_v$.

To see this, assume that to the contrary, for each move $X_v \in \mathcal{M}_t^v$, there is a move $X'_v \in \mathcal{M}_t^v$ such that $X_v \prec_t X'_v$. Then one can construct a sequence of moves $X_v^1, X_v^2, X_v^3, \dots$ such that for all n , $X_v^n \in \mathcal{M}_t^v$ and $X_v^n \prec_t X_v^{n+1}$. By transitivity this implies $X_v^m \prec_t X_v^n$ for all $n > m$. Let $k < \infty$ be the cardinality of \mathcal{M}_t^v . Because of $X_v^g = X_v^h$ for some $g < h \leq k+1$, one finds $X_v^g \prec_t X_v^g$, which contradicts irreflexivity and thus proves that there is some move $X_v \in \mathcal{M}_t^v$ such that for no $Y_v \in \mathcal{M}_t^v$, $X_v \prec_t Y_v$. By lemma 3.1, this must be a t -rational move.

Thus, if node v is reached in some play in Σ_t , there is some play in Σ_t where a t -rational move is played at v . One can now construct a t -rational play $\hat{\sigma}$ in Σ_t by induction on the structure of the game tree as follows. As the root is reached in all plays, there must be a play in Σ_t where a t -rational move at the root is played. Take this as the first move of $\hat{\sigma}$. Now assume you have chosen the first n moves of $\hat{\sigma}$ which are all t -rational and played at some play in Σ_t . Let X_v be the n -th move of $\hat{\sigma}$. If X_v leads to a terminal node, you are finished. Otherwise, let w be the decision node X_v leads to. As X_v is played in some state in Σ_t , clearly $\mathcal{M}_t^w \neq \emptyset$. Thus there is some t -rational move $X_w \in \mathcal{M}_t^w$. Taking X_w as the $n+1$ -th move of $\hat{\sigma}$ completes the induction.

As $\hat{\sigma}$ is t -rational, it belongs to Σ_{t+1} , which is consequently not empty. \blacksquare

Proof (of theorem 3.1). Note that with $\epsilon = 0$, the inequality in the statement of lemma 3.1 simplifies to $\pi_t^{Min}(Y_v) > \pi_t^{Max}(X_v)$.

Let me first show by induction on t that if there is a play in Σ_t where node v is reached, there is also a play in Σ_t containing the backward induction move B_v . For $t = 1$, this is trivial.

Now assume the induction claim has been shown for $t = n$. Let $\pi^{BI}(X_v)$ stand for the payoff of the player playing X_v if in the subgame starting at the node X_v leads to, only backward induction moves are played. If $[X_v]_n \neq \emptyset$, then by the induction hypothesis there is a play $\sigma \in \Sigma_n$ containing X_v and only backward induction moves afterwards. This means $\pi^{BI}(X_v) = \pi_{i(v)}(\sigma) \in \{\pi_{i(v)}(\hat{\sigma}) \mid \hat{\sigma} \in [X_v]_n\}$. With $Y_v \in \mathcal{M}_t^v$, $[Y_v]_n \neq \emptyset$, standing for

any move other than the backward induction move B_v , one therefore finds $\pi_n^{Min}(Y_v) \leq \pi^{BI}(Y_v) < \pi^{BI}(B_v) \leq \pi_n^{Max}(B_v)$, which by lemma 3.1 means that B_v is n -rational.

Now if there is a play in $\sigma \in \Sigma_{n+1}$ where v is reached, this is an n -rational play going through v . Thus one can construct an n -rational play σ' by taking the moves up to v from σ , adding B_v and thereafter choosing moves as in the proof of lemma 2.1. Consequently, there is a state in Σ_{n+1} where B_v is played, which completes the induction.

Define the length $l(v)$ of a decision node v as the maximum number of moves that are played after v has been reached in any play of the game. I will show by induction on $l(v)$ that in any play $\sigma \in \Sigma_t$ at nodes v with $l(v) < t$, only backward induction moves can be played. The base case $t = 1$ is again trivial. For the induction step, assume the induction claim has been shown for all nodes w with $l(w) < t$ and let $l(v) = t$. By the result from the preceding paragraph, if there is a play $\sigma \in \Sigma_t$ where v is reached, there must also be a play $\sigma' \in \Sigma_t$ where the backward induction move B_v is played. For any other move X_v belonging to some play $\sigma \in \Sigma_t$, the induction hypothesis yields $\pi^{BI}(B_v) = \pi_t^{Min}(B_v) > \pi_t^{Max}(X_v) = \pi^{BI}(X_v)$. Therefore by lemma 3.1, X_v is t -irrational and thus does not belong to any play in Σ_{t+1} . ■

Proof (of theorem 3.2). The example before the statement of theorem 3.2 suffices as a proof. ■

REFERENCES

1. Dekel, E., Gul, F.: Rationality and common knowledge in game theory. In: Kreps and Wallis (eds.), *Advances in Economics and Econometrics: Theory and Application*, Vol. I, Cambridge (1996), 87–172
2. Kreps, D., Milgrom, P., Roberts, J., Wilson, R.: Rational Cooperation in the Finitely Repeated Prisoner's Dilemma. *Journal of Economic Theory* **27** (1982) 245–252
3. Lewis, D., *Counterfactuals*, Oxford (1973)
4. McKelvey, R., Palfrey, T.: An Experimental Study of the Centipede Game. *Econometrica* **60** (1992) 803–836
5. Zauner, K.: A Payoff Uncertainty Explanation of Results in Experimental Centipede Games. *Games and Economic Behavior* **26** (1999) 157–185