

This is a section of [doi:10.7551/mitpress/11252.001.0001](https://doi.org/10.7551/mitpress/11252.001.0001)

The Handbook of Rationality

Edited by: Markus Knauff, Wolfgang Spohn

Citation:

The Handbook of Rationality

Edited by: Markus Knauff, Wolfgang Spohn

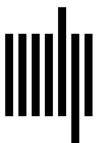
DOI: 10.7551/mitpress/11252.001.0001

ISBN (electronic): 9780262366175

Publisher: The MIT Press

Published: 2021

Funding for the open access edition was provided by the MIT Libraries Open Monograph Fund.



The MIT Press

9.2 Epistemic Game Theory

Andrés Perea

Summary

In this chapter, we review some of the most important ideas, concepts, and results in epistemic game theory, with a focus on the central idea of *common belief in rationality*. We start by showing how belief hierarchies can be encoded by means of epistemic models with types and how this encoding can be used to formally define common belief in rationality. We next indicate how the induced choices can be characterized by a recursive elimination procedure and how the concept relates to Nash equilibrium. Finally, we investigate how the idea of common belief in rationality can be extended to dynamic games by looking at several plausible ways in which players may revise their beliefs.

1. From Classical to Epistemic Game Theory

Classical game theory, as explored in chapter 9.1 by Albert and Kliemt (in this handbook), was pioneered by the seminal work of von Neumann (1928/1959), von Neumann and Morgenstern (1944), and Nash (1950, 1951). It presents a series of concepts that select, for every game and each of the players in that game, a set of possible choices.

In *epistemic game theory*, we concentrate on the *beliefs* that motivate these choices. These may be beliefs about the possible choices of the opponents (*first-order* beliefs) but also beliefs about the *beliefs* of others (*higher-order* beliefs). Putting these first-order and higher-order beliefs together leads to *belief hierarchies*—the language of epistemic game theory. The aim of epistemic game theory is to impose reasonable conditions on such belief hierarchies and to explore the *behavioral consequences* resulting from these conditions.

As with many disciplines in science, it is difficult to say when epistemic game theory really started off. Morgenstern (1935/1976), more than 80 years ago, already stressed the importance of belief hierarchies in economic analysis, but it took a long time before belief hierarchies structurally entered the analysis of human behavior in

economic systems and games. A possible reason for this long delay lies in the complexity of a belief hierarchy. Despite being a very natural object, it is quite difficult to work with because it involves *infinitely many layers*.

The purpose of this chapter is to provide an overview of some of the most important ideas and results in epistemic game theory, with a focus on the central concept of *common belief in rationality*. The outline is as follows: in section 2, we show how infinite belief hierarchies in static games can conveniently be encoded by means of epistemic models with types and use it in section 3 to formally define common belief in rationality. In section 4, we present a recursive elimination procedure that characterizes the choices that can be rationally made under common belief in rationality. In section 5, we discuss the epistemic gap between common belief in rationality and the famous notion of Nash equilibrium. In section 6, finally, we discuss how the idea of common belief in rationality can be extended to dynamic games.

For a more comprehensive overview of epistemic game theory, the reader is referred to the overview paper by Brandenburger (2007), the textbook by Perea (2012), the handbook chapter of Dekel and Siniscalchi (2015), the encyclopedia entry by Pacuit and Roy (2015), and the book by Battigalli, Friedenberg, and Siniscalchi (in press).

2. Belief Hierarchies and Types

The central idea in epistemic game theory is that of *common belief in rationality*. Informally, it states that you do not only choose rationally yourself but also believe that your opponents will choose rationally, that your opponents believe that the other players will choose rationally, and so on. Most other reasoning concepts in epistemic game theory may be viewed as refinements, or variants, of common belief in rationality. The intuitive idea of common belief in rationality is already present in Spohn (1982) and in the concept of rationalizability (Bernheim, 1984; Pearce, 1984), although the latter two papers do not formally define the notion.

An important question is how the idea of common belief in rationality can be defined formally. Consider a *finite static game*¹ $G = (C_i, u_i)_{i \in I}$, where I is the finite set of players, C_i the finite set of choices for player i , and $u_i: \times_{j \in I} C_j \rightarrow \mathbb{R}$ player i 's utility function. It is assumed that all these ingredients—the set of players, the sets of choices, and the utility functions—are *commonly believed* among the players. Moreover, here and in the rest of this chapter, we restrict to noncooperative games. When we say that player i believes in the opponents' rationality, we mean that player i believes that every opponent j chooses optimally, given what player i believes that player j believes about his opponents' choices. For this to be formally defined, we need to specify i 's belief about j 's choice—a *first-order* belief—together with i 's belief about j 's belief about his opponents' choices, which is a *second-order* belief. Similarly, to formally define that player i believes that player j believes in his opponents' rationality, we need to additionally specify the belief that i holds about the belief that j holds about the belief that every opponent k holds about the other players' choices, which is a *third-order* belief. By continuing in this fashion, we can form, for any given player i , an infinite string originating in a first-order belief about the opponents' choices, a second-order belief about the opponents' first-order beliefs, a third-order belief about the opponents' second-order beliefs, and so on. Such infinite strings of beliefs are called *belief hierarchies*. They constitute the central concept of the language of epistemic game theory.

In view of the fact that belief hierarchies are *infinite* strings, making it hard to write them down explicitly, epistemic game theorists typically encode them in an easy and compact way as *types* in the sense of Harsanyi (1967–1968). The main idea is as follows: in a belief hierarchy, player i holds, for every opponent j , a belief about j 's choice, j 's first-order belief, j 's second-order belief, and so on. That is, a belief hierarchy for player i specifies, for every opponent j , a belief about j 's choice and j 's belief hierarchy. If we replace the words “belief hierarchy” by “type” and formalize beliefs by probability distributions, we obtain the following definition of an epistemic model with types:

Definition 1 (Epistemic Model with Types). Consider a finite static game $G = (C_i, u_i)_{i \in I}$. A *finite epistemic model* for G is a tuple $M = (T_i, b_i)_{i \in I}$, where T_i is the finite set of types for player i , and $b_i: T_i \rightarrow \Delta(C_{-i} \times T_{-i})$ is i 's belief mapping, which assigns to every type $t_i \in T_i$ a probabilistic belief $b_i(t_i) \in \Delta(C_{-i} \times T_{-i})$ on the choice–type combinations of i 's opponents.

In this definition, we have used the following pieces of notation: for every finite set X , we denote by $\Delta(X)$ the set

of probability distributions on X . By $C_{-i} \times T_{-i} := \times_{j \neq i} (C_j \times T_j)$, we denote the set of choice–type combinations for i 's opponents.

A finite epistemic model may be viewed as a convenient way to encode belief hierarchies in a finite manner. Indeed, for every type in an epistemic model, we may derive the full belief hierarchy it induces.

To see how this works, consider the two-player game in table 9.2.1, where player 1's choices are in the rows and player 2's choices are in the columns, together with an epistemic model in table 9.2.2.

The superscript of a type always specifies the choice that is optimal for that particular type. This will be shown later. The expression $b_1(t_1^e) = (0.6) \cdot (e, t_2^e) + (0.4) \cdot (f, t_2^e)$ means that type t_1^e assigns probability 0.6 to the event that player 2 chooses e and is of type t_2^e , as well as assigns probability 0.4 to the event that player 2 chooses f and is of type t_2^e . On the other hand, $b_1(t_1^g) = (g, t_2^g)$ means that type t_1^g assigns probability 1 to the event that player 2 chooses g and is of type t_2^g .

Consider the type t_1^b . As t_1^b believes that, with probability 1, player 2 chooses e and is of type t_1^e , the induced first-order belief is that player 1 believes that, with probability 1, player 2 chooses e . Moreover, as player 2's type t_2^e has the belief $(0.6) \cdot (c, t_1^c) + (0.4) \cdot (d, t_1^d)$ about player 1, player 2's type t_2^e assigns probability 0.6 to player 1 choosing c and probability 0.4 to player 1 choosing d . Hence, the second-order belief induced by type t_1^b is that player 1 assigns probability 1 to the event that player 2 assigns probability 0.6 to player 1 choosing c and probability 0.4 to player 1 choosing d . In a similar fashion,

Table 9.2.1
A two-player game

	e	f	g	h
a	0,0	4,1	4,4	4,3
b	3,2	0,0	3,4	3,3
c	2,2	2,1	0,0	2,3
d	1,2	1,1	1,4	0,0

Table 9.2.2
An epistemic model for the game in table 9.2.1

Types	$T_1 = \{t_1^a, t_1^b, t_1^c\}, T_2 = \{t_2^e, t_2^f, t_2^g\}$
Beliefs for player 1	$b_1(t_1^a) = (g, t_2^g)$
	$b_1(t_1^b) = (e, t_2^e)$
	$b_1(t_1^c) = (0.6) \cdot (e, t_2^e) + (0.4) \cdot (f, t_2^f)$
Beliefs for player 2	$b_2(t_2^e) = (0.6) \cdot (c, t_1^c) + (0.4) \cdot (d, t_1^d)$
	$b_2(t_2^f) = (a, t_1^a)$
	$b_2(t_2^g) = (c, t_1^c)$

we can derive the higher-order beliefs, and hence the full belief hierarchy, for the type t_1^b and for all the other types in the epistemic model.

In the game-theoretic literature, people often use *infinite* instead of finite epistemic models, because they wish to work with models that encode *all possible* belief hierarchies. Such exhaustive models are also called *terminal* type structures. That terminal type structures exist for every finite static game—something that is far from obvious—has been shown by Armbruster and Böge (1979), Böge and Eisele (1979), Mertens and Zamir (1985), and Brandenburger and Dekel (1993). For this chapter, we have chosen to work with finite epistemic models instead for two reasons. First, finite epistemic models are easier to work with than terminal type structures, since no advanced measure-theoretic or topological machinery is needed. Moreover, as we will see, this choice does not affect the main results we discuss.

The game-theoretic literature also uses alternative ways of encoding belief hierarchies, such as Kripke structures (Kripke, 1963) and Aumann structures (Aumann, 1974, 1976). The first is the predominant model in the logical and philosophical literature, whereas the latter is often used by economists. Both models use *states* instead of types and assign to every state and every player i a choice for player i , together with a belief for player i about the states. In a similar way as above, one can then derive from such a structure a belief hierarchy for every player at every state. In this chapter, we have chosen to encode belief hierarchies by means of types, but the whole chapter could have been written by using Kripke structures or Aumann structures instead.

3. Common Belief in Rationality

In the previous section, we have seen that belief hierarchies can be encoded by means of epistemic models with types. This now enables us to provide a formal definition of common belief in rationality, starting from the first layer of common belief in rationality, which states that player i believes that every opponent chooses rationally.

To express this event within the formalism of epistemic models with types, we must first define what it means for a choice to be optimal for a type. Consider an epistemic model $M = (T_i, b_i)_{i \in I}$ for a static game $G = (C_i, u_i)_{i \in I}$, a type $t_i \in T_i$, and a choice $c_i \in C_i$. Then,

$$u_i(c_i, t_i) := \sum_{(c_{-i}, t_{-i}) \in C_{-i} \times T_{-i}} b_i(t_i)(c_{-i}, t_{-i}) \cdot u_i(c_i, c_{-i})$$

denotes the *expected utility* for type t_i of choosing c_i . We say that choice c_i is *optimal* for type t_i if $u_i(c_i, t_i) \geq u_i(c'_i, t_i)$ for all $c'_i \in C_i$. In the epistemic model of table 9.2.2, it can

be verified that a is optimal for the type t_1^a , b is optimal for the type t_1^b , and c is optimal for the type t_1^c . Similarly, e is optimal for player 2's type t_2^e , g is optimal for the type t_2^g , and h is optimal for the type t_2^h .

Remember that a type t_i holds a probabilistic belief $b_i(t_i)$ on the opponents' choice–type combinations. For a type t_i to believe in the opponents' rationality means that $b_i(t_i)$ must only assign positive probability to opponents' choice–type pairs where the choice is optimal for the type.

Definition 2 (Belief in the Opponents' Rationality). Consider a finite epistemic model $M = (T_i, b_i)_{i \in I}$ for a finite static game $G = (C_i, u_i)_{i \in I}$. A type $t_i \in T_i$ *believes in the opponents' rationality* if $b_i(t_i)((c_j, t_j)_{j \neq i}) > 0$ only if, for every opponent $j \neq i$, choice c_j is optimal for type t_j .

In the epistemic model of table 9.2.2, it can be verified that types t_1^a , t_1^b , t_2^g , and t_2^h believe in the opponents' rationality, but the other two types do not. Indeed, the type t_1^c for player 1 assigns positive probability to player 2's choice–type pair (f, t_2^e) , where f is not optimal for the type t_2^e , and hence t_1^c does not believe in player 2's rationality. Similarly, player 2's type t_2^e assigns positive probability to player 1's choice–type pair (d, t_1^a) , where d is not optimal for t_1^a , and hence t_2^e does not believe in player 1's rationality.

With the definition of belief in the opponents' rationality at hand, we can now recursively define k -fold belief in rationality for all $k \geq 1$, which finally enables us to formalize common belief in rationality.

Definition 3 (Common Belief in Rationality). Consider a finite epistemic model $M = (T_i, b_i)_{i \in I}$ for a finite static game $G = (C_i, u_i)_{i \in I}$.

(Induction start) A type $t_i \in T_i$ expresses onefold belief in rationality if t_i believes in the opponents' rationality.

(Induction step) For $k > 1$, a type $t_i \in T_i$ expresses k -fold belief in rationality if $b_i(t_i)((c_j, t_j)_{j \neq i}) > 0$ only if, for every opponent $j \neq i$, type t_j expresses $(k - 1)$ -fold belief in rationality.

A type $t_i \in T_i$ expresses *common belief in rationality* if t_i expresses k -fold belief in rationality for every $k \geq 1$.

Hence, twofold belief in rationality entails that a type only assigns positive probability to opponents' types that express onefold belief in rationality. In other words, the player believes that every opponent believes in his opponents' rationality. Similarly, threefold belief in rationality corresponds to the event that a player believes that his opponents believe that their opponents believe in their opponents' rationality, and so on.

Within a finite static game $G = (C_i, u_i)_{i \in I}$, we say that player i can *rationaly choose* $c_i \in C_i$ under common belief in rationality if there is a finite epistemic model $M = (T_i, b_i)_{i \in I}$ and a type $t_i \in T_i$ such that t_i expresses common belief in rationality and c_i is optimal for t_i . That is, choice c_i can be supported by some belief hierarchy that expresses common belief in rationality.

In the epistemic model of table 9.2.2, it can be verified that types t_1^c and t_2^c do not express onefold belief in rationality. Indeed, type t_1^c assigns positive probability to the choice–type pair (f, t_2^g) , where f is not optimal for the type t_2^g and similarly for type t_2^c . Next, types t_1^b and t_2^b express onefold but not twofold belief in rationality. To see why, note that type t_1^b believes that player 2 is of type t_2^c , which does not express onefold belief in rationality, and similarly for type t_2^b . Finally, types t_1^a and t_2^a express common belief in rationality. Indeed, type t_1^a believes that player 2 is of type t_2^g , whereas type t_2^a believes that player 1 is of type t_1^a . As both t_1^a and t_2^a express onefold belief in rationality, it can inductively be shown that t_1^a expresses k -fold belief in rationality for all k and hence expresses common belief in rationality, and similarly for type t_2^a . Consequently, player 1 can rationally choose a and player 2 can rationally choose g under common belief in rationality.

4. Recursive Procedure

Suppose that in a given static game, we are interested in the choices that the players can rationally make under common belief in rationality. Is there an easy method to find these choices, without having to resort to epistemic models with types? That is the question that will be addressed in this section.

The key to answering this question is lemma 3 in Pearce (1984), which we will reproduce below. To state the lemma formally, we need the following definitions. Consider a finite static game $G = (C_i, u_i)_{i \in I}$, a choice c_i , and a belief $b_i \in \Delta(C_{-i})$ about the opponents' choices. Then,

$$u_i(c_i, b_i) := \sum_{c_{-i} \in C_{-i}} b_i(c_{-i}) \cdot u_i(c_i, c_{-i})$$

denotes the expected utility of choice c_i under the belief b_i . Choice c_i is called *optimal* in G for the belief b_i if $u_i(c_i, b_i) \geq u_i(c'_i, b_i)$ for all $c'_i \in C_i$. Choice c_i is called *strictly dominated* in G if there is some randomization $r_i \in \Delta(C_i)$ such that

$$u_i(c_i, c_{-i}) < \sum_{c'_i \in C_i} r_i(c'_i) \cdot u_i(c'_i, c_{-i}) \text{ for all } c_{-i} \in C_{-i}.$$

In the literature, such randomizations $r_i \in \Delta(C_i)$ are typically called *mixed strategies* or *randomized choices*, and they are often interpreted as real objects of choice for player i .

In this chapter, however, we assume that players do not randomize when making a decision, and these randomizations r_i are merely used as an auxiliary device to characterize choices that are optimal for some belief. The reason is that players are assumed to be expected utility maximizers, and hence a player can never increase his expected utility by randomizing over his choices.

Lemma 3 in Pearce (1984) can now be stated as follows:

Lemma 1 (Pearce, 1984). Consider a finite static game $G = (C_i, u_i)_{i \in I}$ and a choice $c_i \in C_i$. Then, there is a belief $b_i \in \Delta(C_{-i})$ such that c_i is optimal in G for b_i if, and only if, c_i is not strictly dominated in G .

This lemma can be used to characterize the choices a player can rationally make if he believes in his opponents' rationality. Let G^1 be the reduced game that remains if we eliminate, for every player, the choices that are strictly dominated in G . For a player to believe in the opponents' rationality thus means, by lemma 1, that his belief is fully concentrated on opponents' choices in G^1 . By applying lemma 1 to the reduced game G^1 , we thus conclude that, for every player, the choices he can rationally make if he believes in the opponents' rationality are exactly the choices in G^1 that are not strictly dominated in G^1 . That is, these are the choices that survive *two rounds* of elimination of strictly dominated choices. In a similar vein, it can be shown that the choices that can rationally be made if a player believes in his opponents' rationality, and believes that his opponents believe in their opponents' rationality (i.e., if he expresses up to *twofold* belief in rationality), are exactly the choices that survive *three rounds* of elimination of strictly dominated choices. By continuing in this fashion, we arrive at the following elimination procedure:

Definition 4 (Iterated Elimination of Strictly Dominated Choices). Consider a finite static game $G = (C_i, u_i)_{i \in I}$.

(Induction start) Let $G^0 := G$ be the full game.

(Induction step) For every $k > 0$, let G^k be the reduced game that remains if we eliminate from G^{k-1} all choices that are strictly dominated in G^{k-1} .

A choice $c_i \in C_i$ survives iterated elimination of strictly dominated choices if c_i is in G^k for all $k > 0$.

By the argument above, we thus see that G^2 contains exactly those choices that can rationally be made if a player believes in the opponents' rationality. By iterating this argument, we conclude that, for every $k \geq 2$, the k -fold reduced game G^k contains exactly those choices that can rationally be made under some belief hierarchy that expresses up to $(k - 1)$ -fold belief in rationality. This

argument already appears in Spohn (1982). In particular, the choices that survive the full procedure will be exactly those choices that can rationally be made under common belief in rationality. This leads to the following central result, which is based on theorems 5.2 and 5.3 in Tan and Werlang (1988), and which Brandenburger (2014) has called the “fundamental theorem of epistemic game theory.” Brandenburger and Dekel (1987) offer in their proposition 2.1 a similar result, characterizing common belief in rationality by “best reply sets” instead of an elimination procedure.

Theorem 1 (Fundamental Theorem of Epistemic Game Theory). Consider a finite static game $G = (C_i, u_i)_{i \in I}$ and a choice $c_i \in C_i$. Then, c_i can rationally be made under common belief in rationality if, and only if, c_i survives iterated elimination of strictly dominated choices.

The fundamental theorem would remain unaffected if we were to use terminal type structures (hence, *infinite* epistemic models) instead of *finite* epistemic models to define common belief in rationality. To illustrate the procedure of iterated elimination of strictly dominated choices and the theorem above, consider the game G from table 9.2.1. In the full game G , it is easily verified that player 1’s choice d is strictly dominated by the randomization that assigns probability 0.5 to his choices a and b and that player 2’s choice f is strictly dominated by the randomization that assigns probability 0.5 to his choices g and h . No other choices are strictly dominated. Hence, G^1 is the game obtained after eliminating the choices d and f . Within the onefold reduced game G^1 , player 1’s choice c is strictly dominated by b (or, rather, the randomization that assigns probability 1 to b), and player 2’s choice e is strictly dominated by h . Hence, G^2 is the game obtained from G^1 after eliminating the choices c and e . Finally, within G^2 , player 1’s choice b is strictly dominated by a , and player 2’s choice h is strictly dominated by g . As such, only the choices a and g survive iterated elimination of strictly dominated choices, and hence, by theorem 1, these are the only choices that can rationally be made under common belief in rationality.

5. Nash Equilibrium

For many decades, the concept of Nash equilibrium (Nash, 1950, 1951) has dominated the classical approach to game theory, inspiring many refinements such as perfect equilibrium (Selten, 1975) and proper equilibrium (Myerson, 1978) for static games, as well as subgame-perfect equilibrium (Selten, 1965) and sequential equilibrium (Kreps

& Wilson, 1982) for dynamic games. See chapter 9.1 by Albert and Kliemt (this handbook) for a discussion of the latter two concepts. However, for a long time, it remained unclear what epistemic conditions are needed for players to choose in accordance with Nash equilibrium. The purpose of this section is to investigate Nash equilibrium from an epistemic perspective and to link it to the conditions of common belief in rationality that we have explored so far. Let us start by giving the definition of Nash equilibrium. See also chapter 9.1, where it is called *mixed-strategy equilibrium*.

Definition 5 (Nash Equilibrium). Consider a finite static game $G = (C_i, u_i)_{i \in I}$. A Nash equilibrium in G is a tuple $(\sigma_i)_{i \in I}$ where $\sigma_i \in \Delta(C_i)$ for every player i , such that $\sigma_i(c_i) > 0$ only if

$$\sum_{c_{-i}=(c_j)_{j \neq i} \in C_{-i}} \left(\prod_{j \neq i} \sigma_j(c_j) \right) \cdot u_i(c_i, c_{-i}) \geq \sum_{c_{-i}=(c'_j)_{j \neq i} \in C_{-i}} \left(\prod_{j \neq i} \sigma_j(c'_j) \right) \cdot u_i(c'_i, c_{-i})$$

for all $c'_i \in C_i$.

In other words, a Nash equilibrium is a tuple of probability distributions on choices such that a choice only receives positive probability if it is optimal against the probability distributions on the opponents’ choices. Traditionally, these probability distributions σ_i have been interpreted as conscious randomizations, or *mixed strategies*, by the players. A more recent approach, adopted by Spohn (1982), Aumann and Brandenburger (1995), and other authors, is to interpret σ_i as the (common) probabilistic belief that i ’s opponents have about i ’s choice, and this is also the interpretation we use here.

A Nash equilibrium $(\sigma_i)_{i \in I}$ induces, in a natural way, a belief hierarchy for player i in which his (first-order) belief about the opponents’ choices is given by $(\sigma_j)_{j \neq i}$, his (second-order) belief about j ’s belief about his opponents’ choices is given by $(\sigma_k)_{k \neq j}$, and so on. Such belief hierarchies are called *simple* in Perea (2012). Moreover, this belief hierarchy can be shown to express common belief in rationality, relying on the optimality conditions in a Nash equilibrium. To see this, consider the belief hierarchy for player i induced by a Nash equilibrium $(\sigma_i)_{i \in I}$. Then, player i only assigns positive probability to a choice c_j of player j if $\sigma_j(c_j) > 0$. By the optimality condition of Nash equilibrium, this is only the case if c_j is optimal against $(\sigma_k)_{k \neq j}$, which is what player i believes that player j believes about his opponents’ choices. Altogether, we see that player i only assigns positive probability to c_j if c_j is optimal for player j , given what player i believes that player j believes about his opponents’ choices. That is, with this belief hierarchy,

player i believes in j 's rationality. In a similar vein, it can be shown that with this belief hierarchy, induced by a Nash equilibrium, player i also believes that every opponent j believes in his opponents' rationality, and so on. Hence, every Nash equilibrium induces, for every player, a belief hierarchy that expresses common belief in rationality. We can thus say that Nash equilibrium implies common belief in rationality.

But is the other direction also true? Does common belief in rationality necessarily lead to Nash equilibrium? The answer, as we will see, is *no*. Consider the two-player game in table 9.2.3, which is similar to the game of figure 1 in Bernheim (1984).

It may be verified that all three choices can rationally be made under common belief in rationality. However, there is only one Nash equilibrium (σ_1, σ_2) in this game, where σ_1 assigns probability 1 to c and σ_2 assigns probability 1 to f . Hence, in this example, Nash equilibrium imposes more restrictions than just common belief in rationality. But what are these extra restrictions?

To see this most clearly, consider the epistemic model, together with its graphical representation, in figure 9.2.1.

It may be verified that all types in the epistemic model express common belief in rationality. Moreover, the respective superscripts of the types indicate the choice that is optimal for that type. Remember that only the choices c and f are supported by a Nash equilibrium in this game.

Consider the type t_i^a that supports the choice a —a choice that is not supported by a Nash equilibrium. The induced belief hierarchy states that, on the one hand, player 1 believes that player 2 chooses e but, on the other hand, believes that 2 believes that 1 believes that 2 chooses d . That is, player 1 believes that player 2 is *incorrect* about 1's first-order belief. The same can be said about his type t_i^b . In contrast, the type t_i^c that supports the Nash equilibrium choice c induces a belief hierarchy in which 1 believes that 2 is correct about 1's first-order belief.

It turns out that in two-player games, this *correct beliefs assumption*—that is, that a player believes that his opponent is correct about his first-order belief—is exactly what separates common belief in rationality

from Nash equilibrium. This is reflected in Spohn's (1982) theorem on page 253, and Aumann and Brandenburger's (1995) theorem A, which both state that in two-player games, mutual belief in rationality, together with mutual belief in the actual first-order beliefs, leads to Nash equilibrium. Here, mutual belief in rationality means that player 1 believes in 2's rationality, and player 2 believes in 1's rationality. Similarly, mutual belief in the actual first-order beliefs means that player 1 is correct about 2's first-order belief, and player 2 is correct about player 1's first-order belief. From a one-person perspective (in which conditions are imposed on the belief hierarchy of a *single* player i) the Spohn–Aumann–Brandenburger conditions thus state that player i believes that j is rational, believes that j believes that i is rational, that i believes that j is correct about i 's first-order belief, and that i believes that j believes that i is correct about j 's first-order belief. In particular, Spohn, Aumann, and Brandenburger show that the *first two layers* of common belief in rationality, together with the correct beliefs assumptions above, are enough to imply Nash equilibrium. Not all layers of common belief in rationality are needed. Polak (1999) shows, however, that if mutual belief in the actual first-order beliefs is strengthened to *common* belief in the actual first-order beliefs, then the Spohn–Aumann–Brandenburger conditions would imply common belief in rationality. Other epistemic foundations for Nash equilibrium in two-player games, which in some way or another involve the correct beliefs assumption above, can be found in Tan and Werlang (1988), Brandenburger and Dekel (1989), Asheim (2006), and Perea (2007a). As the reasonability of the correct beliefs assumption can be debated—after all, why should an opponent be correct about your first-order belief?—these papers implicitly point at the problematic assumptions underlying Nash equilibrium.

For more than two players, the above conditions are no longer enough to characterize Nash equilibrium. For such games, Nash equilibrium additionally implies that i 's belief about j 's choice must be stochastically independent from i 's belief about k 's choice and that i 's belief about j 's belief about k 's choice must be the same as i 's belief about k 's choice. The first property follows from the fact that in a Nash equilibrium $(\sigma_i)_{i \in I}$, the belief of i about the opponents' choices is given by the independent probability distributions $(\sigma_j)_{j \neq i}$, whereas the second condition is implied by the property that i 's belief about j 's belief about k 's choice and i 's belief about k 's choice are both given by σ_k . These two conditions are not implied by common belief in rationality, and hence the gap between Nash equilibrium and common belief in rationality is

Table 9.2.3

A two-player game

	d	e	f
a	0,3	3,0	0,2
b	3,0	0,3	0,2
c	2,0	2,0	2,2

Types	$T_1 = \{t_1^a, t_1^b, t_1^c\}, T_2 = \{t_2^d, t_2^e, t_2^f\}$
Beliefs for player 1	$b_1(t_1^a) = (e, t_2^e)$ $b_1(t_1^b) = (d, t_2^d)$ $b_1(t_1^c) = (f, t_2^f)$
Beliefs for player 2	$b_2(t_2^d) = (a, t_1^a)$ $b_2(t_2^e) = (b, t_1^b)$ $b_2(t_2^f) = (c, t_1^c)$

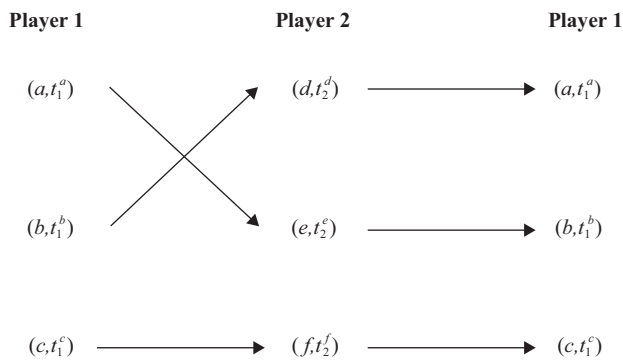


Figure 9.2.1 Epistemic model for the game in table 9.2.3 and a graphical representation.

even bigger in games with more than two players. Epistemic foundations for Nash equilibrium in games with more than two players can be found in Brandenburger and Dekel (1987), Aumann and Brandenburger (1995), Perea (2007a), Barelli (2009), and Bach and Tsakas (2014).

6. Dynamic Games

So far, we have been exploring *static* games, where all players only make one choice, and players choose in complete ignorance of the opponents' choices. We now investigate how the idea of common belief in rationality can be translated to *dynamic games*. In a dynamic game, players may choose one after the other, may choose more than once, and may fully or partially observe what the opponents have done in the past when it is their turn to move. As a consequence, a player may need to *revise* his belief about the opponents when he discovers that his previous belief has been contradicted by some of the opponents' past choices. As an illustration, consider the game from figure 9.2.2, which is based on Reny (1992).

If player 1 believes that player 2 would rationally choose g at his last move, then he would choose a at the beginning. Common belief in rationality thus seems to suggest that player 2 should initially believe that player 1 chooses a . However, when it is player 2's turn to move, this initial belief has been contradicted by player 1's past play, and hence player 2 must revise his belief about player 1. But how? As we will see, there are at least two plausible ways for player 2 to revise his belief.

One option is to interpret player 1's past move b as a *mistake*, yet at the same time maintain the belief that player 1 would choose rationally at his second move, as well as the belief that player 1 believes that player 2 would rationally choose g at his second move. In that case, player 2 would believe, upon observing b , that player 1 would choose e at this second move, and therefore player 2 would choose c . This type of reasoning, in which the players are free to interpret "surprising" past moves as mistakes but believe that the opponents will choose rationally in the future, believe that the opponents always believe that their opponents will choose rationally in the future, and so on, is called *backward induction reasoning* and is formally captured by the concept of *common belief in future rationality* (Perea, 2014). Similar lines of reasoning are present in Penta (2015), Baltag, Smets, and Zvesper (2009), and the concept of *sequential rationalizability* (Dekel, Fudenberg, & Levine, 1999, 2002; Asheim & Perea, 2005). Backward induction reasoning is also implicitly present in the backward induction procedure (for a survey of the various epistemic foundations for backward induction, see Perea, 2007b) and the equilibrium concepts of subgame-perfect equilibrium (Selten, 1965) and sequential equilibrium (Kreps & Wilson, 1982; for a formal statement, see Perea & Predtetchinski, 2019).

Another option for player 2, after observing the "surprising" move b , is to interpret b as a conscious, optimal choice for player 1. However, this is only possible if player 2 believes that player 1 would choose f afterward and if player 2 believes that player 1 assigns a high probability to player 2 making the suboptimal choice h at his second move. Consequently, player 2 would choose d and, in case he is asked to make a second move, choose g . This type of reasoning, where a player, whenever possible, tries to interpret "surprising" past choices as conscious, optimal choices, is called *forward induction reasoning*. It can be formalized by the epistemic condition of *strong belief in the opponents' rationality* (Battigalli & Siniscalchi, 2002), which states that a player, whenever possible, must believe that his opponents are implementing optimal strategies.² The concepts that most closely implement this type of reasoning are *explicable equilibrium*

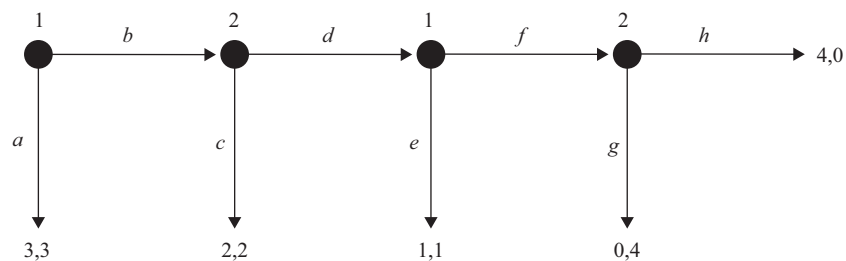


Figure 9.2
Reny's game.

(Reny, 1992) and *extensive-form rationalizability* (Pearce, 1984; Battigalli, 1997), where the latter has epistemically been characterized by *common strong belief in rationality* in Battigalli and Siniscalchi (2002). See also Battigalli and Friedenberg (2012), in which forward induction with exogenous restrictions on the players' beliefs is studied.

As the example above illustrates, backward and forward induction reasoning may lead to different strategy choices. Indeed, player 2 chooses c under backward induction reasoning but would choose (d, g) under forward induction reasoning. However, both types of reasoning lead to the same outcome, which is the terminal history following a . Battigalli (1997) has shown that the latter is always true in dynamic games with perfect information without relevant ties by proving that in every such game, the forward induction concept of extensive-form rationalizability always uniquely leads to the backward induction outcome. This result is remarkable, as forward induction and backward induction represent two completely different lines of reasoning. The connection between these two lines of reasoning in general dynamic games is one of the many intriguing problems in epistemic game theory that need further exploration.

Notes

1. Such a static game may also correspond to the strategic form of a dynamic game, as explained in chapter 9.1 by Albert and Kliemt (this handbook).
2. Strong belief in rationality is very similar to *assumption of rationality* in static games, which has been used by Brandenburger, Friedenberg, and Keisler (2008) to epistemically characterize the iterated elimination of weakly dominated choices.

References

- Armbruster, W., & Böge, W. (1979). Bayesian game theory. In O. Moeschlin & D. Pallaschke (Eds.), *Game theory and related topics* (pp. 17–28). Amsterdam, Netherlands: North-Holland.
- Asheim, G. B. (2006). *The consistent preferences approach to deductive reasoning in games* (Theory and Decision Library). Dordrecht, Netherlands: Springer.

Asheim, G. B., & Perea, A. (2005). Sequential and quasi-perfect rationalizability in extensive games. *Games and Economic Behavior*, 53, 15–42.

Aumann, R. J. (1974). Subjectivity and correlation in randomized strategies. *Journal of Mathematical Economics*, 1, 67–96.

Aumann, R. J. (1976). Agreeing to disagree. *Annals of Statistics*, 4, 1236–1239.

Aumann, R., & Brandenburger, A. (1995). Epistemic conditions for Nash equilibrium. *Econometrica*, 63, 1161–1180.

Bach, C. W., & Tsakas, E. (2014). Pairwise epistemic conditions for Nash equilibrium. *Games and Economic Behavior*, 85, 48–59.

Baltag, A., Smets, S., & Zvesper, J. A. (2009). Keep “hoping” for rationality: A solution to the backward induction paradox. *Synthese*, 169, 301–333.

Barelli, P. (2009). Consistency of beliefs and epistemic conditions for Nash and correlated equilibrium. *Games and Economic Behavior*, 67, 363–375.

Battigalli, P. (1997). On rationalizability in extensive games. *Journal of Economic Theory*, 74, 40–61.

Battigalli, P., & Friedenberg, A. (2012). Forward induction reasoning revisited. *Theoretical Economics*, 7, 57–98.

Battigalli, P., Friedenberg, A., & Siniscalchi, M. (in press). *Epistemic game theory: Reasoning about strategic uncertainty*.

Battigalli, P., & Siniscalchi, M. (2002). Strong belief and forward induction reasoning. *Journal of Economic Theory*, 106, 356–391.

Bernheim, B. D. (1984). Rationalizable strategic behavior. *Econometrica*, 52, 1007–1028.

Böge, W., & Eisele, T. (1979). On solutions of Bayesian games. *International Journal of Game Theory*, 8, 193–215.

Brandenburger, A. (2007). The power of paradox: Some recent developments in interactive epistemology. *International Journal of Game Theory*, 35, 465–492.

Brandenburger, A. (Ed.). (2014). *The language of game theory: Putting epistemics into the mathematics of games* (World Scientific Series in Economic Theory, Vol. 5). Singapore: World Scientific.

Brandenburger, A., & Dekel, E. (1987). Rationalizability and correlated equilibria. *Econometrica*, 55, 1391–1402.

- Brandenburger, A., & Dekel, E. (1989). The role of common knowledge assumptions in game theory. In F. Hahn (Ed.), *The economics of missing markets, information and games* (pp. 46–61). Oxford, England: Oxford University Press.
- Brandenburger, A., & Dekel, E. (1993). Hierarchies of beliefs and common knowledge. *Journal of Economic Theory*, *59*, 189–198.
- Brandenburger, A., Friedenberg, A., & Keisler, H. J. (2008). Admissibility in games. *Econometrica*, *76*, 307–352.
- Dekel, E., Fudenberg, D., & Levine, D. K. (1999). Payoff information and self-confirming equilibrium. *Journal of Economic Theory*, *89*, 165–185.
- Dekel, E., Fudenberg, D., & Levine, D. K. (2002). Subjective uncertainty over behavior strategies: A correction. *Journal of Economic Theory*, *104*, 473–478.
- Dekel, E., & Siniscalchi, M. (2015). Epistemic game theory. In P. Young & S. Zamir (Eds.), *Handbook of game theory with economic applications* (Vol. 4, pp. 619–702). Amsterdam, Netherlands: North-Holland.
- Harsanyi, J. C. (1967–1968). Games with incomplete information played by “Bayesian” players, I–III, *Management Science*, *14*, 159–182, 320–334, 486–502.
- Kreps, D. M., & Wilson, R. (1982). Sequential equilibria. *Econometrica*, *50*, 863–894.
- Kripke, S. (1963). A semantical analysis of modal logic I: Normal modal propositional calculi. *Zeitschrift für mathematische Logik und Grundlagen der Mathematik*, *9*, 67–96.
- Mertens, J.-F., & Zamir, S. (1985). Formulation of Bayesian analysis for games with incomplete information. *International Journal of Game Theory*, *14*, 1–29.
- Morgenstern, O. (1976). Perfect foresight and economic equilibrium. In A. Schotter (Ed.), *Selected economic writings of Oskar Morgenstern* (pp. 169–183). New York: New York University Press. (Original work published 1935)
- Myerson, R. B. (1978). Refinements of the Nash equilibrium concept. *International Journal of Game Theory*, *7*, 73–80.
- Nash, J. F., Jr. (1950). Equilibrium points in n -person games. *Proceedings of the National Academy of Sciences of the United States of America*, *36*, 48–49.
- Nash, J. F. (1951). Non-cooperative games. *Annals of Mathematics*, *54*, 286–295.
- Pacuit, E., & Roy, O. (2015). Epistemic foundations of game theory. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*. Retrieved from <https://plato.stanford.edu/archives/spr2015/entries/epistemic-game/>
- Pearce, D. G. (1984). Rationalizable strategic behavior and the problem of perfection. *Econometrica*, *52*, 1029–1050.
- Penta, A. (2015). Robust dynamic implementation. *Journal of Economic Theory*, *160*, 280–316.
- Perea, A. (2007a). A one-person doxastic characterization of Nash strategies. *Synthese*, *158*, 251–271.
- Perea, A. (2007b). Epistemic foundations for backward induction: An overview. In J. van Benthem, D. Gabbay, & B. Löwe (Eds.), *Interactive logic: Proceedings of the 7th Augustus de Morgan Workshop, London* (Texts in Logic and Games, Vol. 1, pp. 159–193). Amsterdam, Netherlands: Amsterdam University Press.
- Perea, A. (2012). *Epistemic game theory: Reasoning and choice*. Cambridge, England: Cambridge University Press.
- Perea, A. (2014). Belief in the opponents’ future rationality. *Games and Economic Behavior*, *83*, 231–254.
- Perea, A., & Predtetchinski, A. (2019). An epistemic approach to stochastic games. *International Journal of Game Theory*, *48*, 181–203.
- Polak, B. (1999). Epistemic conditions for Nash equilibrium, and common knowledge of rationality. *Econometrica*, *67*, 673–676.
- Reny, P. J. (1992). Backward induction, normal form perfection and explicable equilibria. *Econometrica*, *60*, 627–649.
- Selten, R. (1965). Spieltheoretische Behandlung eines Oligopolmodells mit Nachfragezeit [Game-theoretic treatment of an oligopoly model with demand time]. *Zeitschrift für die gesamte Staatswissenschaft*, *121*, 301–324, 667–689.
- Selten, R. (1975). Reexamination of the perfectness concept for equilibrium points in extensive games. *International Journal of Game Theory*, *4*, 25–55.
- Spohn, W. (1982). How to make sense of game theory. In W. Stegmüller, W. Balzer, & W. Spohn (Eds.), *Philosophy of economics: Proceedings, Munich, July 1981* (pp. 239–270). Berlin, Germany: Springer.
- Tan, T. C.-C., & Werlang, S. R. da C. (1988). The Bayesian foundations of solution concepts of games. *Journal of Economic Theory*, *45*, 370–391.
- von Neumann, J. (1959). Zur Theorie der Gesellschaftsspiele. *Mathematische Annalen*, *100*, 295–320. Translated by S. Bargmann as “On the theory of games of strategy (S. Bargmann, Trans.). In A. W. Tucker & R. D. Luce (Eds.), *Contributions to the theory of games* (Vol. IV, pp. 13–43). Princeton, NJ: Princeton University Press. (Original work published 1928)
- von Neumann, J., & Morgenstern, O. (1944). *Theory of games and economic behavior*. Princeton, NJ: Princeton University Press.

© 2021 The Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

The MIT Press would like to thank the anonymous peer reviewers who provided comments on drafts of this book. The generous work of academic experts is essential for establishing the authority and quality of our publications. We acknowledge with gratitude the contributions of these otherwise uncredited readers.

This book was set in Stone Serif and Stone Sans by Westchester Publishing Services.

Library of Congress Cataloging-in-Publication Data

Names: Knauff, Markus, editor. | Spohn, Wolfgang, editor.

Title: The handbook of rationality / edited by Markus Knauff and Wolfgang Spohn.

Description: Cambridge : The MIT Press, 2021. | Includes bibliographical references and index.

Identifiers: LCCN 2020048455 | ISBN 9780262045070 (hardcover)

Subjects: LCSH: Reasoning (Psychology) | Reason. | Cognitive psychology. | Logic. | Philosophy of mind.

Classification: LCC BF442 .H36 2021 | DDC 153.4/3—dc23

LC record available at <https://lcn.loc.gov/2020048455>