

Bonanno, Giacomo

**Working Paper**

## A dynamic epistemic characterization of backward induction without counterfactuals

Working Paper, No. 12-2

**Provided in Cooperation with:**

University of California Davis, Department of Economics

*Suggested Citation:* Bonanno, Giacomo (2012) : A dynamic epistemic characterization of backward induction without counterfactuals, Working Paper, No. 12-2, University of California, Department of Economics, Davis, CA

This Version is available at:

<https://hdl.handle.net/10419/58360>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

Giacomo Bonanno  
UC Davis

March 18, 2012

Paper # 12-2

The analysis of rational play in dynamic games is usually done within a static framework that specifies a player's initial beliefs as well as his disposition to revise those beliefs conditional on hypothetical states of information. We suggest a simpler approach, where the rationality of a player's choice is judged on the basis of the actual beliefs that the player has at the time he has to make that choice. We propose a dynamic framework where the set of "possible worlds" is given by state-instant pairs  $(w, t)$ . Each state  $w$  specifies the entire play of the game and, for every instant  $t$ ,  $(w, t)$  specifies the history that is reached at that instant (in state  $w$ ). A player is said to be active at  $(w, t)$  if the history reached in state  $w$  at date  $t$  is a decision history of his. At every state-instant pair  $(w, t)$  the beliefs of the active player provide an answer to the question "what will happen if I take action  $a$ ", for every available action  $a$ . A player is said to be rational at  $(w, t)$  if either he is not active there or the action he ends up taking at state  $w$  is "optimal" given his beliefs at  $(w, t)$ . We provide a characterization of backward induction in terms of the following event: the first mover (i) is rational and has correct beliefs, (ii) believes that the active player at date 1 is rational and has correct beliefs, (iii) believes that the active player at date 1 believes that the active player at date 2 is rational and has correct beliefs, etc. Thus our epistemic characterization does not rely on dispositional belief revision or on (objective or subjective) counterfactuals.

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

[http://www.econ.ucdavis.edu/working\\_search.cfm](http://www.econ.ucdavis.edu/working_search.cfm)

# A dynamic epistemic characterization of backward induction without counterfactuals

Giacomo Bonanno

Department of Economics,  
University of California,  
Davis, CA 95616-8578 - USA  
gfbonanno@ucdavis.edu

March 2012

## Abstract

The analysis of rational play in dynamic games is usually done within a static framework that specifies a player's initial beliefs as well as his disposition to revise those beliefs conditional on hypothetical states of information. We suggest a simpler approach, where the rationality of a player's choice is judged on the basis of the actual beliefs that the player has at the time he has to make that choice. We propose a dynamic framework where the set of "possible worlds" is given by state-instant pairs  $(\omega, t)$ . Each state  $\omega$  specifies the entire play of the game and, for every instant  $t$ ,  $(\omega, t)$  specifies the history that is reached at that instant (in state  $\omega$ ). A player is said to be active at  $(\omega, t)$  if the history reached in state  $\omega$  at date  $t$  is a decision history of his. At every state-instant pair  $(\omega, t)$  the beliefs of the active player provide an answer to the question "what will happen if I take action  $a$ ?", for every available action  $a$ . A player is said to be rational at  $(\omega, t)$  if either he is not active there or the action he ends up taking at state  $\omega$  is optimal given his beliefs at  $(\omega, t)$ . We provide a characterization of backward induction in terms of the following event: the first mover (i) is rational and has correct beliefs, (ii) believes that the active player at date 1 is rational and has correct beliefs, (iii) believes that the active player at date 1 believes that the active player at date 2 is rational and has correct beliefs, etc. Thus our epistemic characterization does not rely on dispositional belief revision or on (objective or subjective) counterfactuals.

Keywords: perfect-information game, backward induction, dynamic interactive beliefs, rationality, Kripke frame

# 1 Introduction

The analysis of rational play in dynamic games is usually done within a *static* framework that specifies, for every player, his initial beliefs as well as his disposition to revise those beliefs conditional on hypothetical states of information that the player might find himself in. This is done by means of interactive structures which model a rather complex web of beliefs: for example, Player 2 might initially believe that Player 1 will end the game right away and yet have very detailed beliefs about what Player 1 would believe about Player 2's revised beliefs if Player 1 were instead to give the move to Player 2. In these models each player is assumed to have not only a disposition to revise his own beliefs, should he be faced with unexpected information, but also to have (conditional) beliefs about the disposition of the other players to revise their beliefs. This seems to constitute a rather "heavy weight" approach to modeling the players' states of mind in a dynamic game. It is shown in this literature ([6, 8, 13, 15]) that common *initial* belief of rationality does not imply a backward induction outcome in perfect-information games.

In this paper we suggest an alternative and simpler approach, where the rationality of a player's choice is judged on the basis of the *actual beliefs* that the player has *at the time he has to make that choice*. We propose a dynamic analysis of perfect-information games where the set of "possible worlds" is given by state-instant pairs  $(\omega, t)$ . Each state  $\omega$  specifies the entire play of the game and, for every instant  $t$ ,  $(\omega, t)$  specifies the history that is reached at that instant (in state  $\omega$ ). A player is said to be *active* at  $(\omega, t)$  if the history reached in state  $\omega$  at date  $t$  is a decision history of his. At every state-instant pair  $(\omega, t)$  the beliefs of the active player provide an answer to the question "what will happen if I take action  $a$ ?", for every available action  $a$ . A player is said to be rational at  $(\omega, t)$  if either he is not active there or the action he ends up taking at state  $\omega$  is "optimal" given his beliefs at  $(\omega, t)$ . We provide a characterization of backward induction in terms of the following event: the first mover (i) is rational and has correct beliefs, (ii) believes that the active player at date 1 is rational and has correct beliefs, (iii) believes that the active player at date 1 believes that the active player at date 2 is rational and has correct beliefs, etc.

This can be stated more precisely as follows. First we define a time- $t$  belief operator  $B_t$  which captures the beliefs of the active player and enables us to express a player's belief that the next player will respond rationally to his choice. Let  $\mathbf{T}_t$  be the set of states where the active player at date  $t$  (if there is any) has correct beliefs and let  $\mathbf{R}_t$  be the set of states where the choice of the active player at date  $t$  is rational. In keeping with the literature, we focus on perfect-information games with no relevant ties where there is a unique backward-induction solution. We prove the following characterization. If  $\omega \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0(\mathbf{T}_1 \cap \mathbf{R}_1) \cap B_0 B_1(\mathbf{T}_2 \cap \mathbf{R}_2) \cap \dots \cap B_0 B_1 \dots B_{m-2}(\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$  (where  $m$  is the depth of the game) then the play associated with  $\omega$  is the backward-induction play. Conversely, if  $z$  is the backward-induction play then there is a model and a state  $\omega$  such that  $\omega \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0(\mathbf{T}_1 \cap \mathbf{R}_1) \cap \dots \cap B_0 B_1 \dots B_{m-2}(\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$  and the play associated with  $\omega$  is  $z$ . Thus

we obtain an *epistemic characterization of backward induction that does not rely on (objective or subjective) counterfactuals or on dispositional belief revision*.

## 2 Perfect-information games and models

We use the history-based definition of extensive-form game. If  $A$  is a set, we denote by  $A^*$  the set of finite sequences in  $A$ . If  $h = \langle a_1, \dots, a_k \rangle \in A^*$  and  $1 \leq j \leq k$ , the sequence  $\langle a_1, \dots, a_j \rangle$  is called a *prefix* of  $h$ . If  $h = \langle a_1, \dots, a_k \rangle \in A^*$  and  $a \in A$ , we denote the sequence  $\langle a_1, \dots, a_k, a \rangle \in A^*$  by  $ha$ .

A *finite extensive form with perfect information* (without chance moves) is a tuple  $\langle A, H, N, \iota, \rangle$  whose elements are:

- A finite set of actions  $A$ .
- A finite set of histories  $H \subseteq A^*$  which is closed under prefixes (that is, if  $h \in H$  and  $h' \in A^*$  is a prefix of  $h$ , then  $h' \in H$ ). The null history  $\langle \rangle$ , denoted by  $\emptyset$ , is an element of  $H$  and is a prefix of every history. A history  $h \in H$  such that, for every  $a \in A$ ,  $ha \notin H$ , is called a *terminal history*. The set of terminal histories is denoted by  $Z$ .  $D = H \setminus Z$  denotes the set of non-terminal or *decision* histories. For every history  $h \in D$ , we denote by  $A(h)$  the set of actions available at  $h$ , that is,  $A(h) = \{a \in A : ha \in H\}$ .
- A finite set  $N = \{1, \dots, n\}$  of players.
- A function  $\iota : D \rightarrow N$  that assigns a player to each decision history. Thus  $\iota(h)$  is the player who moves at history  $h$ . For every  $i \in N$ , let  $D_i = \iota^{-1}(i)$  be the set of histories assigned to player  $i$ .

Given an extensive form, one obtains an *extensive game* by adding, for every player  $i \in N$ , a *utility* (or *payoff*) *function*  $U_i : Z \rightarrow \mathbb{R}$  (where  $\mathbb{R}$  denotes the set of real numbers; recall that  $Z$  is the set of terminal histories).

Given a history  $h \in H$ , we denote by  $\ell(h)$  the length of  $h$ , which is defined recursively as follows:  $\ell(\emptyset) = 0$  and if  $h \in D$  and  $a \in A(h)$  then  $\ell(ha) = \ell(h) + 1$ . Thus  $\ell(h)$  is equal to the number of actions that appear in  $h$ ; for example, if  $h = \langle \emptyset, a_1, a_2, a_3 \rangle$  then  $\ell(h) = 3$ . We denote by  $\ell^{\max}$  the length of the maximal history in  $H$ :  $\ell^{\max} = \max_{h \in H} \ell(h)$ . Clearly, if  $\ell(h) = \ell^{\max}$  then  $h \in Z$  (that is,  $\max_{h \in H} \ell(h) = \max_{h \in Z} \ell(h)$ ). Given a history  $h \in H$  and an integer  $t$  with  $0 \leq t \leq \ell^{\max}$ , we denote by  $h_t$  the prefix of  $h$  of length  $t$ . For example, if  $h = \langle \emptyset, a, b, c, d \rangle$ , then  $h_0 = \emptyset$ ,  $h_2 = \langle \emptyset, a, b \rangle$ , etc.

From now on histories will be denoted more succinctly by listing the corresponding actions, without angled brackets and without commas: thus instead of writing  $\langle \emptyset, a_1, a_2, a_3, a_4 \rangle$  we will simply write  $a_1 a_2 a_3 a_4$ .

Let  $\Omega$  be a set of states and  $T$  a set of instants or dates. We call the set  $\Omega \times T$  the set of *state-instant pairs*. If  $E \subseteq \Omega \times T$  and  $t \in T$ , we denote by  $E_t$  the set of states  $\{\omega \in \Omega : (\omega, t) \in E\}$ .

**Definition 1** Given an extensive form with perfect information  $G = \langle A, H, N, \iota, \rangle$ , a state-time representation of  $G$  is a triple  $\langle \Omega, T, \zeta \rangle$  where  $\Omega$  is a set of states,  $T = \{0, 1, \dots, \ell^{\max}\}$  and  $\zeta : \Omega \rightarrow Z$  is a function that assigns to every state a terminal history. Given a state-instant pair  $(\omega, t) \in \Omega \times T$ , let

$$\zeta_t(\omega) = \begin{cases} \text{the prefix of } \zeta(\omega) \text{ of length } t & \text{if } t < \ell(\zeta(\omega)) \\ \zeta(\omega) & \text{if } t \geq \ell(\zeta(\omega)). \end{cases}$$

Interpretation: the play of the game unfolds over time; the first move is made at date 0, the second move at date 1, etc. A state  $\omega \in \Omega$  specifies a particular play of the game (that is, a complete sequence of moves leading to terminal history  $\zeta(\omega)$ );  $\zeta_t(\omega)$  denotes the “state of play at time  $t$ ” in state  $\omega$ , that is, the partial history of the play up to date  $t$  [if  $t$  is less than the length of  $\zeta(\omega)$ , otherwise - once the play is completed - the state of the system remains at  $\zeta(\omega)$ ]. Figure 1 shows an extensive form with perfect information and a state-time representation of it. For every  $\omega \in \Omega = \{\alpha, \beta, \gamma\}$  and  $t \in T = \{0, 1, 2, 3\}$  we have indicated the (partial) history  $\zeta_t(\omega)$  (recall that  $\emptyset$  denotes the empty history).

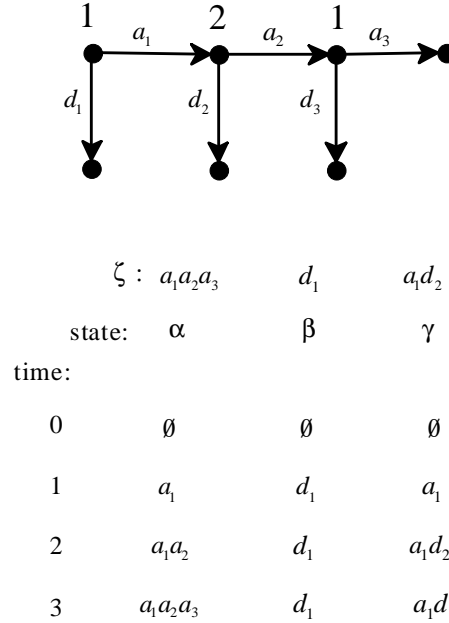


Figure 1: An extensive form with perfect information and a state-time representation of it

We want to define the notion of rational behavior in a game and examine its implications. Player  $i$  chooses rationally at a decision history of his, if the choice he makes there is “optimal” given the beliefs that he holds *at the time*

at which he makes that choice. These beliefs might be different from his initial beliefs about what would happen in the game and thus might be revised beliefs in light of the information he has at the moment. However, his prior beliefs are *not relevant* in assessing the rationality of his choice: what counts is what he believes at the time he makes the decision. The beliefs (prior or revised) of the other players are also irrelevant. Thus in order to assess the rationality of the actual behavior of the players all we need to specify at a state-instant pair  $(\omega, t)$  are the *actual* beliefs of the *active* player. This can be done within a state-time representation of the game. Given a state  $\omega$  and an instant  $t$ , there will be a unique player who makes a decision at  $(\omega, t)$  (unless the play of the game has already reached a terminal history, in which case there are no decisions to be made). If  $\zeta_t(\omega)$  is a decision history, the active player is  $\iota(\zeta_t(\omega))$ ; denote  $\zeta_t(\omega)$  by  $h$  and  $\iota(\zeta_t(\omega))$  by  $i$ . Then player  $i$  has to choose an action from the set  $A(h)$ . In order to make this choice he will form some beliefs about what will happen if he chooses action  $a$ , for every  $a \in A(h)$ . These beliefs will be used to assess the rationality of the choice that the player ends up making at state  $\omega$ . We will describe a player's beliefs about the consequences of taking alternative actions by means of an accessibility relation. Thus we use Kripke frames and represent qualitative, rather than probabilistic, beliefs. In order to simplify the notation, we will assign beliefs also to the non-active players, but in a trivial way by making those players believe everything.

We recall the following facts about Kripke frames. If  $\Omega$  is a set of states and  $\mathcal{B}_i \subseteq \Omega \times \Omega$  a binary relation on  $\Omega$  (representing the beliefs of individual  $i$ ), for every  $\omega \in \Omega$  we denote by  $\mathcal{B}_i(\omega)$  the set of states that are reachable from  $\omega$  using  $\mathcal{B}_i$ , that is,  $\mathcal{B}_i(\omega) = \{\omega' \in \Omega : \omega \mathcal{B}_i \omega'\}$ .  $\mathcal{B}_i$  is *serial* if  $\mathcal{B}_i(\omega) \neq \emptyset$ , for every  $\omega \in \Omega$ ; it is *transitive* if  $\omega' \in \mathcal{B}_i(\omega)$  implies  $\mathcal{B}_i(\omega') \subseteq \mathcal{B}_i(\omega)$  and it is *euclidean* if  $\omega' \in \mathcal{B}_i(\omega)$  implies  $\mathcal{B}_i(\omega) \subseteq \mathcal{B}_i(\omega')$ . Subsets of  $\Omega$  are called *events*. If  $E \subseteq \Omega$  is an event, we say that at  $\omega \in \Omega$  individual  $i$  believes  $E$  if and only if  $\mathcal{B}_i(\omega) \subseteq E$ . Thus one can define a *belief operator*  $B_i : 2^\Omega \rightarrow 2^\Omega$  as follows:  $B_i E = \{\omega \in \Omega : \mathcal{B}_i(\omega) \subseteq E\}$ . Hence  $B_i E$  is the event that individual  $i$  believes  $E$ . It is well known that seriality of  $\mathcal{B}_i$  corresponds to consistency of beliefs (if the individual believes  $E$  then it is not the case that he believes not  $E$  :  $B_i E \subseteq \neg B_i \neg E$ , where, for every event  $F$ ,  $\neg F$  denotes the complement of  $F$  in  $\Omega$ ), transitivity corresponds to positive introspection (if the individual believes  $E$  then he believes that he believes  $E$  :  $B_i E \subseteq B_i B_i E$ ) and euclideanness corresponds to negative introspection (if the individual does not believe  $E$  then he believes that he does not believe  $E$  :  $\neg B_i E \subseteq B_i \neg B_i E$ ).<sup>1</sup>

**Definition 2** Given an extensive form with perfect information  $G$ , a model of  $G$  is a tuple  $\langle \Omega, T, \zeta, \{\mathcal{B}_{i,t}\}_{i \in N, t \in T} \rangle$  where  $\langle \Omega, T, \zeta \rangle$  is a state-time representation of  $G$  (see Definition 1) and, for every player  $i \in N$  and instant  $t \in T$ ,  $\mathcal{B}_{i,t} \subseteq \Omega \times \Omega$  is a binary relation on the set of states (representing the beliefs of player  $i$  at time  $t$ ) that satisfies the following properties:  $\forall i \in N, \forall t \in T, \forall \omega, \omega', \omega'' \in \Omega$ ,

---

<sup>1</sup>For more details see [5].

1. If  $i \neq \iota(\zeta_t(\omega))$ , that is, if  $\zeta_t(\omega)$  is not a decision history of player  $i$ , then  $\mathcal{B}_{i,t}(\omega) = \emptyset$ .
2. If  $i = \iota(\zeta_t(\omega))$ , that is, if  $\zeta_t(\omega)$  is a decision history of player  $i$ , then
  - 2.1.  $\mathcal{B}_{i,t}$  is locally serial, transitive and euclidean,<sup>2</sup>
  - 2.2. If  $\omega' \in \mathcal{B}_{i,t}(\omega)$  then  $\zeta_t(\omega') = \zeta_t(\omega)$ ,
  - 2.3. For every  $a \in A(\zeta_t(\omega))$  there exists an  $\omega' \in \mathcal{B}_{i,t}(\omega)$  such that  $\zeta_{t+1}(\omega') = \zeta_t(\omega')a$ .

Condition 1 says that a player has trivial beliefs (that is, he believes everything) at all the state-instant pairs where he is not active. We impose this condition only for notational convenience, to eliminate the need to keep track, at every state-instant pair, of who the active player is.<sup>3</sup> To understand Condition 2, fix a state-instant pair  $(\omega, t)$ , let  $h = \zeta_t(\omega)$  and suppose that  $h$  is a decision history of player  $i$  (thus  $i = \iota(\zeta_t(\omega))$ ) where he has to choose an action from the set  $A(h)$ . Condition 2.1 says that player  $i$  has beliefs with standard properties (consistency, positive and negative introspection). Condition 2.2 says that every state  $\omega'$  which is accessible from  $\omega$  by  $\mathcal{B}_{i,t}$  (that is, every state that player  $i$  considers possible) is such that the history associated with  $(\omega', t)$  is still  $h$ ; in other words, player  $i$  at time  $t$  knows that his decision history  $h$  has been reached. Condition 2.3 says that for every action  $a$  available at  $h$ , there is a state  $\omega'$  that player  $i$  considers possible ( $\omega' \in \mathcal{B}_{i,t}(\omega)$ ) where he takes action  $a$ , that is, the truncation of  $\zeta(\omega')$  at time  $t+1$  (namely  $\zeta_{t+1}(\omega')$ ) is equal to  $ha$  (recall that, by Condition 2.2,  $\zeta_t(\omega') = h$ ). This means that, for every available action, player  $i$  has a belief about what will happen if he chooses that action.

**Remark 3** *It is worth noting that this way of modeling beliefs is a departure from the standard (static) approach in the literature, where it is assumed that if a player takes a particular action at a state then he knows that he takes that action. The standard approach thus requires the use of either objective or subjective counterfactuals in order to represent a player's beliefs about the consequences of taking alternative actions. In our approach a player's beliefs refer to the deliberation or pre-choice stage, where the player considers the consequences of all his actions, without pre-judging his subsequent decision. The state encodes the player's actual choice which can be judged to be rational or irrational by relating it to his pre-choice beliefs. Thus it is possible for a player to have the same beliefs in two different states, say  $\alpha$  and  $\beta$ , and be labeled as rational at state  $\alpha$  and irrational at state  $\beta$ , because the action he ends up taking at state  $\alpha$  is optimal given those beliefs, while the action he ends up taking at state  $\beta$  is not optimal given those same beliefs.*

<sup>2</sup>That is,  $\mathcal{B}_{i,t}(\omega) \neq \emptyset$  and if  $\omega' \in \mathcal{B}_{i,t}(\omega)$  then  $\mathcal{B}_{i,t}(\omega') = \mathcal{B}_{i,t}(\omega)$ .

<sup>3</sup>As explained below, by defining  $\mathcal{B}_t = \bigcup_{i \in N} \mathcal{B}_{i,t}$ , we can take the relation  $\mathcal{B}_t$  to be a description of the beliefs of the active player at date  $t$  (whose identity can change from state to state).



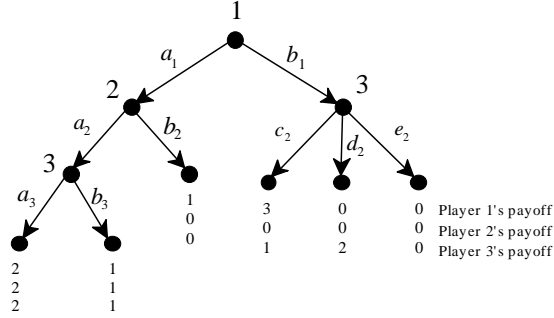


Figure 2: A perfect-information game

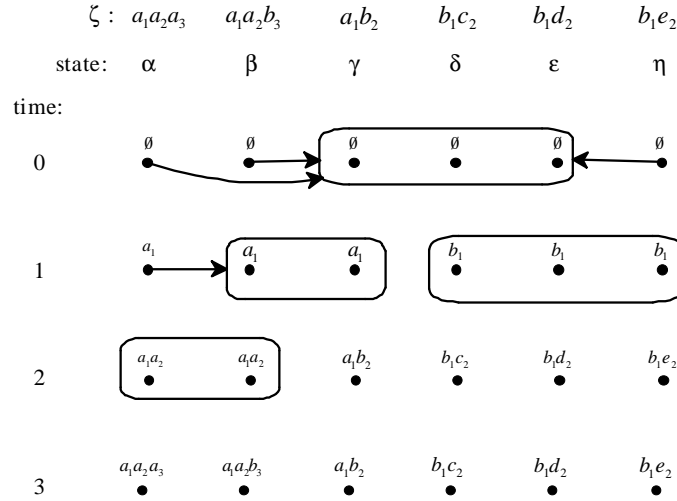


Figure 3: A model of the game of Figure 2

Figure 2 shows a perfect information game and Figure 3 a model of it. We represent a belief relation  $\mathcal{B}$  as follows: for any two states  $\omega$  and  $\omega'$ ,  $\omega' \in \mathcal{B}(\omega)$  if and only if either  $\omega$  and  $\omega'$  are enclosed in the same rounded rectangle or there is an arrow from  $\omega$  to the rounded rectangle containing  $\omega'$ .<sup>4</sup> The relations shown in Figure 3 are those of the active players: the relation at date 0 is that of Player 1 ( $\mathcal{B}_{1,0}$ ), the relation at date 1 for states  $\alpha, \beta$  and  $\gamma$  is that of Player 2 ( $\mathcal{B}_{2,1}$ ), the relation at date 1 for states  $\delta, \varepsilon$  and  $\eta$  is that of Player 3 ( $\mathcal{B}_{3,1}$ ) and the relation at date 2 for states  $\alpha$  and  $\beta$  is that of Player 3 ( $\mathcal{B}_{3,2}$ ).<sup>5</sup> Consider a

<sup>4</sup>In other words, for any two states  $\omega$  and  $\omega'$  that are enclosed in a rounded rectangle,  $\{(\omega, \omega'), (\omega, \omega'), (\omega', \omega), (\omega', \omega')\} \subseteq \mathcal{B}$  (that is, the relation is total on the set of states contained in the rectangle) and if there is an arrow from a state  $\omega$  to a rounded rectangle then, for every  $\omega'$  in the rectangle,  $(\omega, \omega') \in \mathcal{B}$ .

<sup>5</sup>Thus  $\mathcal{B}_{1,0}(\omega) = \{\gamma, \delta, \varepsilon\}$  for every  $\omega \in \Omega$ ,  $\mathcal{B}_{2,1}(\omega) = \{\beta, \gamma\}$  for every  $\omega \in \{\alpha, \beta, \gamma\}$ ,

state, say  $\alpha$ . Then  $\alpha$  describes the following beliefs: at date 0 Player 1 believes that if she takes action  $a_1$  then Player 2 will follow (at date 1) with  $b_2$  (state  $\gamma$ ) and if she takes action  $b_1$  then Player 3 will follow (at date 1) with either  $c_2$  (state  $\delta$ ) or  $d_2$  (state  $\varepsilon$ ); at date 1 Player 2 (knows that Player 1 played  $a_1$  and) believes that if he takes action  $a_2$  then Player 3 will follow (at date 2) with  $b_3$  (and if he takes action  $b_2$  the game will end). At state  $\alpha$  Player 1 ends up playing  $a_1$ , Player 2 ends up playing  $a_2$  and Player 3 ends up playing  $a_3$ .

It is worth noting that the notion of model that we are using allows for erroneous beliefs. Indeed, in the model of Figure 3, at state  $\alpha$  Player 1 has incorrect beliefs about the subsequent move of Player 2 if she herself plays  $a_1$ .

### 3 Rationality and backward induction

We say that at a state-instant pair  $(\omega, t)$  a player is rational if either she is not active at  $\zeta_t(\omega)$  (that is,  $\zeta_t(\omega)$  is not a decision history of hers) or the action that she ends up choosing at  $\omega$  is “optimal” given her beliefs at date  $t$ , in the sense that it is not the case that - according to her beliefs - there is another action of hers that yields higher utility.<sup>6</sup> Thus a player is *irrational* at a state-instant pair  $(\omega, t)$  if she is active at history  $\zeta_t(\omega)$ , she ends up taking action  $a$  at  $\omega$  and she believes that her maximum utility if she takes action  $a$  is less than the minimum utility that she gets if she takes some other action  $a'$ .

**Definition 4** Fix an arbitrary player  $i$  and an arbitrary state-instant pair  $(\omega, t)$ . We say that player  $i$  is rational at  $(\omega, t)$  if either

- (1)  $\zeta_t(\omega)$  is not a decision history of player  $i$ , or
- (2)  $\zeta_t(\omega)$  is a decision history of player  $i$  and if  $a$  is the action chosen by player  $i$  at  $\omega$  (that is,  $\zeta_{t+1}(\omega) = \zeta_t(\omega)a$ ) then, for every  $a' \in A(\zeta_t(\omega))$ , it is not the case that  $\min_{\omega' \in A'} U_i(\zeta(\omega')) > \max_{\omega' \in A} U_i(\zeta(\omega'))$  where  $A' = \{\omega' \in \mathcal{B}_{i,t}(\omega) : \zeta_{t+1}(\omega') = \zeta_t(\omega')a'\}$  and  $A = \{\omega' \in \mathcal{B}_{i,t}(\omega) : \zeta_{t+1}(\omega') = \zeta_t(\omega')a\}$  (recall that  $U_i : Z \rightarrow \mathbb{R}$  is player  $i$ 's utility function on the set of terminal histories).

For example, in the model of Figure 3, Player 1 is rational at state  $\alpha$  and date 0, because she believes that if she takes action  $a_1$  then her payoff will be 1 (she believes that Player 2 will follow with  $b_2$ ) and if she takes action  $b_1$  then her payoff will be either 3 or 0 (she believes that Player 3 will follow with either  $c_2$  or  $d_2$ ) and she actually ends up taking action  $a_1$ . Similarly, Player 2 is rational at state  $\alpha$  and date 1 and Player 3 is rational at state  $\alpha$  and date 2. On the other hand, Player 2 is not rational at state  $\gamma$  and date 1 (he believes that if he takes action  $a_2$  his payoff will be 1 and if he takes action  $b_2$  his payoff will be 0

---

$\mathcal{B}_{3,1}(\omega) = \{\delta, \varepsilon, \eta\}$  for every  $\omega \in \{\delta, \varepsilon, \eta\}$  and  $\mathcal{B}_{3,2}(\omega) = \{\alpha, \beta\}$  for every  $\omega \in \{\alpha, \beta\}$ . For any remaining state  $\omega$  and date  $t$ ,  $\mathcal{B}_{i,t}(\omega) = \emptyset$ , for every player  $i$ . Thus, for example,  $\mathcal{B}_{1,1}(\omega) = \mathcal{B}_{1,2}(\omega) = \mathcal{B}_{1,3}(\omega) = \emptyset$ , for every state  $\omega$ .

<sup>6</sup>This notion of rationality has been referred to in the literature as “material rationality” (see, for example, [2, 6]).

and yet ends up taking action  $b_2$ ). Thus, since  $\gamma \in \mathcal{B}_{1,0}(\alpha)$ , at state  $\alpha$  and date 0 it is not the case that Player 1 believes that Player 2 will choose rationally at date 1.

We denote by  $\mathbf{R}_t \subseteq \Omega$  the event that (i.e. the set of states at which) the active player (if there is one) is rational at date  $t$ .<sup>7</sup> Thus  $\omega \in \mathbf{R}_t$  if and only if either  $\zeta_t(\omega)$  is a terminal history (that is,  $\zeta_t(\omega) = \zeta(\omega)$ ) or  $\zeta_t(\omega)$  is a decision history and the active player at  $\zeta_t(\omega)$  is rational at  $(\omega, t)$ . Of course, the identity of the active player can vary across states, that is, the active player at  $(\omega, t)$  can be different from the active player at  $(\omega', t)$ . In the model of Figure 3 we have that  $\mathbf{R}_0 = \Omega$ ,  $\mathbf{R}_1 = \{\alpha, \beta, \varepsilon\}$ ,  $\mathbf{R}_2 = \{\alpha, \gamma, \delta, \varepsilon, \eta\}$  and  $\mathbf{R}_3 = \Omega$ .

Let  $B_{i,t} : 2^\Omega \rightarrow 2^\Omega$  be the belief operator of player  $i$  at date  $t$ . Thus, for every event  $E \subseteq \Omega$ ,  $B_{i,t}E = \{\omega \in \Omega : \mathcal{B}_{i,t}(\omega) \subseteq E\}$ . By (1) of Definition 2, if player  $i$  is not active at  $(\omega, t)$  then  $\mathcal{B}_{i,t}(\omega) = \emptyset$  and thus  $\omega \in B_{i,t}E$  for every event  $E$ . Let  $B_t : 2^\Omega \rightarrow 2^\Omega$  be the operator defined by  $B_tE = \bigcap_{i \in N} B_{i,t}E$  (thus  $\omega \in B_tE$  if and only if  $\bigcup_{i \in N} \mathcal{B}_{i,t}(\omega) \subseteq E$ ). Then  $B_tE$  is the event that “the active player believes  $E$  at date  $t$ ” (which is trivially equivalent to the event that “everybody believes  $E$  at date  $t$ ”). For example, in the model of Figure 3, we have that  $\alpha \notin B_0\mathbf{R}_1$  (since  $\gamma \in \mathcal{B}_0(\alpha)$  and  $\gamma \notin \mathbf{R}_1$ ), that is, it is not the case that the active player at date 0 (Player 1) believes that the active player at date 1 will choose rationally. Indeed, Player 1 believes that if she plays  $a_1$  then the active player at date 1 (Player 2) will not choose rationally (given the date-1 beliefs that Player 1 ascribes to Player 2) and if she plays  $b_1$  then the active player at date 1 (Player 3) might or might not choose rationally (Player 3 chooses rationally at  $\varepsilon$  but not at  $\delta$ ).

Note that the models we are considering allow for the possibility that a player may ascribe to a future mover beliefs that are different from the beliefs that that player will actually have. In other words, a player may have erroneous beliefs about the future beliefs of other players (or even about her own future beliefs).

Let  $\mathbf{T}_t$  be the set of states where the beliefs of the active player (if there is one) are correct:  $\mathbf{T}_t = \{\omega \in \Omega : \text{if } \mathcal{B}_t(\omega) \neq \emptyset \text{ then } \omega \in \mathcal{B}_t(\omega)\}$ . For example, in the model of Figure 3 we have that  $\mathbf{T}_0 = \{\gamma, \delta, \varepsilon\}$ ,  $\mathbf{T}_1 = \{\beta, \gamma, \delta, \varepsilon, \eta\}$  and  $\mathbf{T}_2 = \mathbf{T}_3 = \Omega$ . Thus if  $\omega \in \mathbf{T}_t$  and  $\zeta_t(\omega)$  is a decision history, then, for every event  $E$ , if the active player believes  $E$  (that is, if  $\mathcal{B}_t(\omega) \subseteq E$ ) then  $E$  is indeed the case (that is,  $\omega \in E$ ).

The following two propositions provide a characterization of backward induction<sup>8</sup> in terms of the following event (where  $m$  is the depth of the game, that is, the length of its maximal histories):

$$(\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0(\mathbf{T}_1 \cap \mathbf{R}_1) \cap B_0B_1(\mathbf{T}_2 \cap \mathbf{R}_2) \cap \dots \cap B_0B_1 \dots B_{m-2}(\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$$

In keeping with the literature, we restrict attention to games without relevant ties.

<sup>7</sup>By Definition 4 inactive players are always rational; thus  $\mathbf{R}_t$  can also be described as the event that “every player is rational at date  $t$ ”.

<sup>8</sup>The definition of backward-induction solution is reviewed in the Appendix.

**Definition 5** A perfect-information game has no relevant ties if,  $\forall i \in N, \forall h \in D_i, \forall a, a' \in A(h)$  with  $a \neq a', \forall z, z' \in Z$ , if  $ha$  is a prefix of  $z$  and  $ha'$  is a prefix of  $z'$  then  $U_i(z) \neq U_i(z')$ .

For example, the game shown in Figure 2 has no relevant ties. If a game has no relevant ties, then it has a unique backward-induction solution.

The proofs of the following propositions are given in the Appendix.

**Proposition 6** Fix a perfect-information game  $G$  without relevant ties and let  $m$  be its depth. Fix an arbitrary model of  $G$  and an arbitrary state  $\omega$ . If  $\omega \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0 (\mathbf{T}_1 \cap \mathbf{R}_1) \cap B_0 B_1 (\mathbf{T}_2 \cap \mathbf{R}_2) \cap \dots \cap B_0 B_1 \dots B_{m-2} (\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$  then  $\zeta(\omega)$  is the backward-induction terminal history.

**Proposition 7** Fix a perfect-information game  $G$  without relevant ties and let  $m$  be its depth. Let  $z$  be the backward-induction terminal history. Then there is a model of  $G$  and a state  $\omega$  such that (1)  $\zeta(\omega) = z$  and (2)  $\omega \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0 (\mathbf{T}_1 \cap \mathbf{R}_1) \cap \dots \cap B_0 B_1 \dots B_{m-2} (\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$ .

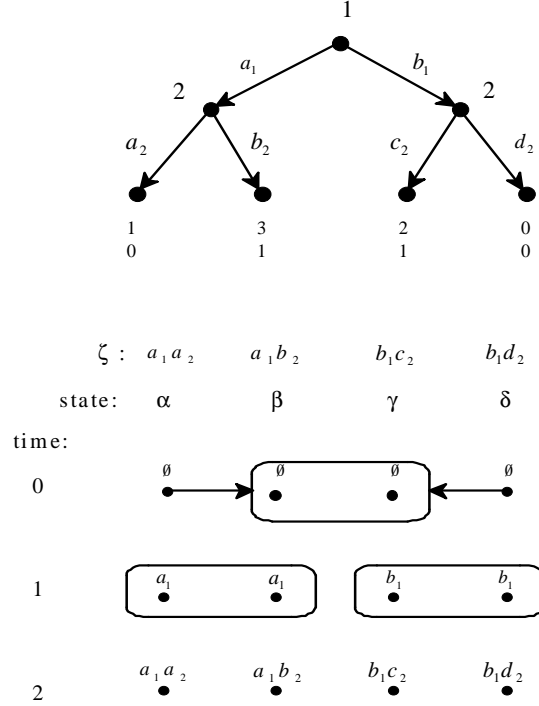


Figure 4: A perfect-information game and a model of it

The condition in Proposition 6 that beliefs be locally correct is essential. For example, if  $\omega \notin \mathbf{T}_0$  then it may happen that  $\zeta(\omega)$  is not the backward-induction terminal history, as shown in Figure 4. Here we have that  $\mathbf{R}_0 = \{\alpha, \beta\}$ ,  $\mathbf{R}_1 = \{\beta, \gamma\}$ ,  $\mathbf{T}_0 = \{\beta, \gamma\}$ ,  $\mathbf{T}_1 = \Omega$ ,  $B_0 \mathbf{R}_1 = B_0 \mathbf{T}_1 = \Omega$ . Hence  $\alpha \in \mathbf{R}_0 \cap B_0 (\mathbf{T}_1 \cap \mathbf{R}_1)$

(in this game  $m = 2$ ) and yet  $\zeta(\alpha) = a_1a_2$  which is not the backward-induction play. At state  $\alpha$  Player 1 is rational, believes that after his move Player 2 will be rational and will have correct beliefs and yet the play associated with  $\alpha$  is not the backward-induction play (because Player 1 is wrong in her belief that Player 2 will play rationally at date 1). Similar examples can be constructed to show that in Proposition 6 the condition  $\omega \in B_0\mathbf{T}_1$  is necessary and so is  $\omega \in B_0B_1\mathbf{T}_2$ , etc.

## 4 Comparison with the literature

There is a large literature on the epistemic foundations of backward induction, which was recently reviewed in [9, 11]. In what follows we shall try to highlight the important differences between our approach and the existing literature.

We have focused on a purely behavioral framework, where a state describes the actual play of the game at the histories that are actually reached in that state. Thus, contrary to a well-established literature ([1, 2, 3, 8, 10, 12, 15]), strategies (or plans of action) do not play any role in our analysis. Indeed the use of strategies in models of dynamic games involves the implicit use of counterfactuals.<sup>9</sup> Methodologically, this is not satisfactory: if it is necessary to specify what a player would do in situations that do not occur in the state under consideration, then one should model the counterfactual explicitly.

The purely behavioral point of view that we have adopted (consisting in associating with every state a play of the game rather than a strategy profile) was first introduced in [13]. Unlike the other papers that take a purely behavioral point of view ([4, 6, 7, 13]), our analysis does not make use of objective or subjective counterfactuals and belief revision plays no role. The use of subjective counterfactuals or dispositional belief revision is made necessary in that literature by two characteristics of the models used. First of all, the static nature of the framework makes it impossible to model explicitly the beliefs of players over time; one thus needs to do so indirectly by representing simultaneously the initial beliefs and the disposition to revise those beliefs subject to conceivable items of information that one might receive during the play of the game. This is done either probabilistically by means of conditional probability systems ([6, 7]) or by means of qualitative belief revision structures ([4]). As pointed out by Stalnaker,

"It should be noted that even with the addition of the belief revision structure to the epistemic models ..., they remain static models. A model of this kind represents only the agent's beliefs at a fixed time, together with the policies or dispositions to revise her beliefs that she has at that time. The model does not represent

---

<sup>9</sup>While in a simultaneous game the association of a strategy of player  $i$  to a state can be interpreted as a description of player  $i$ 's behavior at that state, in the case of dynamic games this interpretation is no longer valid, since one would end up describing not only the actual behavior of player  $i$  but also his counterfactual behavior at a different state (that is, at decision histories that are not reached in the actual state).

any actual revisions that are made when new information is actually received.” ([16], p. 198.)<sup>10</sup>

The second characteristic of the models that use subjective counterfactuals is that they impose the constraint that if at a state a player takes action  $a$  then she knows that she takes action  $a$ ; that is, at every state that the player considers possible, she takes action  $a$ . Thus one needs to use either objective or subjective counterfactuals in order to represent a player’s beliefs about the consequences of taking an action different from  $a$ .

As pointed out in Remark 3, in our approach a player’s beliefs refer to the deliberation or *pre-choice stage*, where the player considers the consequences of all his actions, without pre-judging his subsequent decision. Thus the beliefs of the active player at a state-instant pair are truly open to the possibility of taking any of the available actions: one cannot reason towards a choice if one already knows what that choice will be.

The characterization of backward induction that we have provided is in terms of the forward beliefs of the active player at date 0 (the first mover): she believes in the rationality of future movers and believes that they, too, will believe in the rationality of future movers. That this type of condition is central to backward induction is now well understood ([3, 4, 12, 15]). The novelty of our approach lies in (1) the switch to a dynamic framework for beliefs, (2) showing that the notion of backward induction does not require the use of (objective or subjective) counterfactuals and (3) pointing out the need for “knowledge”, interpreted - locally - as true belief.<sup>11</sup>

## 5 Appendix

We provide below the proofs of Propositions 6 and 7. First we recall the definition of *backward induction solution*. The backward induction solution of a perfect-information game without relevant ties is unique and is given by the output of the following algorithm:

---

<sup>10</sup>The author goes on to say that “The models can be enriched by adding a temporal dimension to represent the dynamics, but doing so requires that the knowledge and belief operators be time indexed...” In our models the belief operators are indeed time indexed and represent the actual beliefs of the players when actually informed that it is their turn to move.

<sup>11</sup>A strand in the literature ([7, 1, 2, 3]) assumes that each belief relation is reflexive everywhere, so that it gives rise to a partition. In such cases it is common to speak of knowledge rather than belief. As Stalanker points out, it is methodologically preferable to carry out the analysis in terms of (possibly erroneous) beliefs and then - if desired - add further conditions, such as the local correctness of beliefs. The reason why one should not start with the assumption of necessarily correct beliefs (that is, global reflexivity of the belief relations) is that such an assumption has strong intersubjective implications:

“The assumption that Alice believes (with probability one) that Bert believes (with probability one) that the cat ate the canary tells us nothing about what Alice believes about the cat and the canary themselves. But if we assume instead that Alice knows that Bert knows that the cat ate the canary, it follows, not only that the cat in fact ate the canary, but that Alice knows it, and therefore believes it as well.” ([14], p. 153.)

1. Start at a decision history  $h$  whose immediate successors are only terminal histories, that is, for every  $a \in A(h)$ ,  $ha \in Z$  (e.g. history  $b_1$  in the game of Figure 2) and select the choice that maximizes the utility of player  $i(h)$  (in the game of Figure 2, at  $b_1$  player 3' utility-maximizing choice is  $d_2$ ). Delete the successors of  $h$ , thus turning  $h$  into a terminal history, and assign to  $h$  the payoff vector associated with the selected choice.
2. Repeat Step 1 in the reduced game until all the decision histories have been exhausted.

The output of the backward-induction algorithm can be written in terms of a profile of strategies, where a strategy of player  $i$  is defined as a list of choices, one for each decision history of player  $i$ . For example, the backward induction solution of the game of Figure 2 can be written as  $(a_1, a_2, (a_3, d_2))$ .

In order to prove Proposition 6 we need the following definition.

**Definition 8** Fix a perfect-information game and a model of it. Let  $\alpha, \beta \in \Omega$ . We say that  $\beta$  is reachable from  $\alpha$  with  $s$  steps ( $s \geq 1$ ) if there is a sequence of state-instant pairs  $\langle (\omega_0, 0), (\omega_1, 1), \dots, (\omega_s, s) \rangle$  such that:

1.  $\omega_0 = \alpha$ ,
2.  $\omega_s = \beta$ ,
3.  $\forall k = 1, \dots, s, \omega_k \in \mathcal{B}_{k-1}(\omega_{k-1})$ .

For example, in the model of Figure 3  $\beta$  is reachable from  $\eta$  with 2 steps with the sequence  $\langle (\eta, 0), (\gamma, 1), (\beta, 2) \rangle$ .<sup>12</sup>

**Remark 9** Let  $E$  be an event,  $\alpha$  a state and suppose that  $\alpha \in B_0 B_1 \dots B_{s-1} E$ . Then for every  $\beta \in \Omega$ , if  $\beta$  is reachable from  $\alpha$  with  $s$  steps then  $\beta \in E$ .<sup>13</sup>

**Proof of Proposition 6.** Fix a perfect-information game with no relevant ties, so that there is a unique backward induction (BI) solution. Let  $f_{BI} : H \rightarrow Z$  be the following function: if  $h$  is a decision history then  $f_{BI}(h)$  is the terminal history that is reached from  $h$  by following the backward-induction choices and if  $z$  is a terminal history then  $f_{BI}(z) = z$ . Recall that if  $z \in Z$  and  $t \in T$ , we denote by  $z_t$  the prefix of  $z$  of length  $t$  (see Definition 1). Fix a model of the game and suppose that  $\alpha$  is a state such that  $\alpha \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0 (\mathbf{T}_1 \cap \mathbf{R}_1) \cap B_0 B_1 (\mathbf{T}_2 \cap \mathbf{R}_2) \cap \dots \cap B_0 B_1 \dots B_{m-2} (\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$  (where  $m$  is the depth of the game). We need to show that  $\zeta(\alpha) = f_{BI}(\emptyset)$  (recall that  $\emptyset$  denotes the empty history). First we show that,

<sup>12</sup>Note that, if  $\beta$  is reachable from  $\alpha$  with  $s$  steps, then  $\zeta_{s-1}(\beta)$  is a decision history. In fact, we have that  $\beta = \omega_s \in \mathcal{B}_{s-1}(\omega_{s-1})$  and thus  $\mathcal{B}_{s-1}(\omega_{s-1}) \neq \emptyset$ , so that  $\zeta_{s-1}(\omega_{s-1})$  is a decision history. (Note also that, by Definition 2,  $\zeta_{s-1}(\beta) = \zeta_{s-1}(\omega_{s-1})$ .)

<sup>13</sup>Proof. Let  $\langle (\omega_0, 0), (\omega_1, 1), \dots, (\omega_s, s) \rangle$  be a sequence that satisfies the properties of Definition 8. Then, since  $\alpha \in B_0 B_1 B_2 \dots B_{s-1} E$ ,  $\mathcal{B}_0(\alpha) \subseteq B_1 B_2 \dots B_{s-1} E$ ; thus, since  $\omega_1 \in \mathcal{B}_0(\alpha)$ ,  $\omega_1 \in B_1 B_2 \dots B_{s-1} E$ . Thus  $\mathcal{B}_1(\omega_1) \subseteq B_2 \dots B_{s-1} E$ , etc. Thus  $\mathcal{B}_{s-1}(\omega_{s-1}) \subseteq E$  and hence, since  $\beta = \omega_s$  and  $\omega_s \in \mathcal{B}_{s-1}(\omega_{s-1})$ ,  $\beta \in E$ .

For every  $t$  with  $1 \leq t \leq m-1$  and for every  $\beta \in \Omega$ ,  
if  $\beta$  is reachable from  $\alpha$  with  $t$  steps then  $\zeta(\beta) = f_{BI}(\zeta_t(\beta))$ . (1)

We prove this by induction.

Base step:  $t = m-1$ . Fix an arbitrary  $\beta$  which is reachable from  $\alpha$  with  $m-1$  steps. If  $\zeta_{m-1}(\beta)$  is a terminal history, then  $\zeta_{m-1}(\beta) = \zeta(\beta)$  (see Definition 1) and, by definition of  $f_{BI}(\cdot)$ ,  $f_{BI}(\zeta(\beta)) = \zeta(\beta)$ . Thus  $\zeta(\beta) = f_{BI}(\zeta_{m-1}(\beta))$ . Suppose, therefore, that  $\zeta_{m-1}(\beta)$  is a decision history. Let  $i$  be the active player, that is, the player who moves at  $\zeta_{m-1}(\beta)$ . Fix an arbitrary  $\omega \in \mathcal{B}_{m-1}(\beta)$  (recall that, by Definition 2,  $\mathcal{B}_{m-1}(\beta) \neq \emptyset$ ). Then, by Definition 2,  $\zeta_{m-1}(\omega) = \zeta_{m-1}(\beta)$ . Since the depth of the game is  $m$ , after player  $i$ 's move at  $\zeta_{m-1}(\omega)$  the game ends and thus  $\zeta_m(\omega) = \zeta(\omega)$ . Since  $\alpha \in B_0 B_1 \dots B_{m-2} \mathbf{R}_{m-1}$ , by Remark 9  $\beta \in \mathbf{R}_{m-1}$ , that is, player  $i$  is rational at state  $\beta$  and time  $m-1$ . Hence, the choice made by player  $i$  at state  $\beta$  and time  $m-1$  is the payoff-maximizing choice there, that is,  $\zeta(\beta) = f_{BI}(\zeta_{m-1}(\beta))$ .

Induction step: suppose that (1) is true for  $t = k$  with  $1 < k \leq m-1$ . We want to show that it is true for  $t = k-1$ . Fix an arbitrary  $\beta$  which is reachable from  $\alpha$  with  $k-1$  steps.

First we show that,

$$\forall \omega \in \mathcal{B}_{k-1}(\beta), \quad \zeta(\omega) = f_{BI}(\zeta_k(\omega)) \quad (2)$$

If  $\zeta_{k-1}(\beta)$  is a terminal history, there is nothing to prove, since  $\mathcal{B}_{k-1}(\beta) = \emptyset$ . Suppose, therefore, that  $\zeta_{k-1}(\beta)$  is a decision history. Let  $i$  be the active player, that is, the player who moves at  $\zeta_{k-1}(\beta)$ . Fix an arbitrary  $\omega \in \mathcal{B}_{k-1}(\beta)$ . Then, by Definition 8,  $\omega$  is reachable from  $\alpha$  with  $k$  steps (since, by hypothesis,  $\beta$  is reachable from  $\alpha$  with  $k-1$  steps). By the induction hypothesis  $\zeta(\omega) = f_{BI}(\zeta_k(\omega))$ . Thus (2) holds. Since  $\alpha \in B_0 B_1 \dots B_{k-2} \mathbf{T}_{k-1}$ , by Remark 9  $\beta \in \mathbf{T}_{k-1}$ ; thus, since  $\mathcal{B}_{k-1}(\beta) \neq \emptyset$ ,

$$\beta \in \mathcal{B}_{k-1}(\beta). \quad (3)$$

Since  $\alpha \in B_0 B_1 \dots B_{k-2} \mathbf{R}_{k-1}$ , by Remark 9  $\beta \in \mathbf{R}_{k-1}$ , that is, player  $i$  is rational at state  $\beta$  and time  $k-1$ . By (2) at state  $\beta$  and time  $k-1$  player  $i$  believes that after his move the play will continue according to the BI solution. Hence the action chosen by  $i$  at  $\zeta_{k-1}(\beta)$  is the optimal action there given those beliefs (i.e. the action dictated by the BI solution), that is,

$$\zeta_k(\beta) \text{ is a prefix of } f_{BI}(\zeta_{k-1}(\beta)). \quad (4)$$

By (2) and (3),  $\zeta(\beta) = f_{BI}(\zeta_k(\beta))$ . It follows from this and (4) that  $\zeta(\beta) = f_{BI}(\zeta_{k-1}(\beta))$ . This completes the proof of (1).

Next we show that

$$\forall \omega \in \mathcal{B}_0(\alpha), \quad \zeta(\omega) = f_{BI}(\zeta_1(\omega)). \quad (5)$$



Fix an arbitrary  $\omega \in \mathcal{B}_0(\alpha)$ . Then  $\omega$  is reachable from  $\alpha$  with 1 step and thus, by (1),  $\zeta(\omega) = f_{BI}(\zeta_1(\omega))$ . Thus the active player at state  $\alpha$  and date 0 believes that after her move the play will continue according to the BI solution. Since  $a \in \mathbf{R}_0$ , it follows that the action chosen by the active player at  $\zeta_0(\alpha) = \emptyset$  is the optimal action there given those beliefs, that is,

$$\zeta_1(\alpha) \text{ is a prefix of } f_{BI}(\emptyset). \quad (6)$$

Since  $\mathcal{B}_0(\alpha) \neq \emptyset$  and  $\alpha \in \mathbf{T}_0$ ,  $\alpha \in \mathcal{B}_0(\alpha)$ . Thus, by (5),  $\zeta(\alpha) = f_{BI}(\zeta_1(\alpha))$ . It follows from this and (6) then  $\zeta(\alpha) = f_{BI}(\emptyset)$ . ■

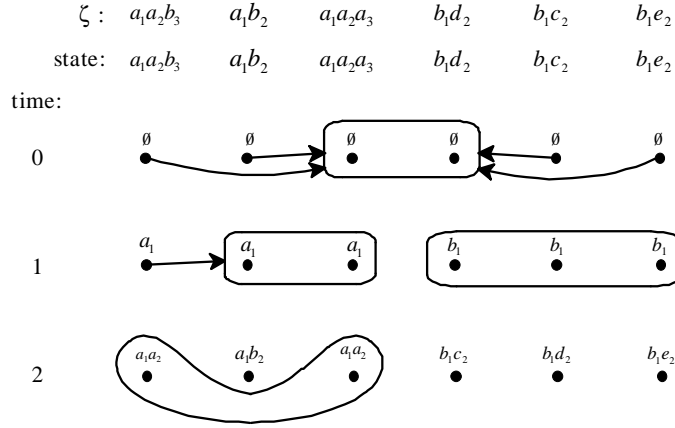


Figure 5

A model of the game of Figure 2.

**Proof of Proposition 7.** Fix a perfect-information game  $G$  and define the following model of it:  $\Omega = Z$  (recall that  $Z$  is the set of terminal histories),  $T = \{0, 1, \dots, m = \ell^{\max}\}$  (recall that  $\ell^{\max}$  is the depth of the game, that is, the length of its maximal histories) and  $\zeta$  is the identity function (that is,  $\zeta(z) = z$ , for every  $z \in Z$ ). Let  $f_{BI} : H \rightarrow Z$  be the function defined in the proof of Proposition 6. Fix an arbitrary player  $i$ , an arbitrary  $z \in Z$  and an arbitrary  $t \in T$ . If  $z_t$  is not a decision history of player  $i$ , then we set  $\mathcal{B}_{i,t}(z) = \emptyset$ ; if  $z_t$  is a decision history of player  $i$  then we set  $\mathcal{B}_{i,t}(z) = \{z' \in Z : z'_t = z_t \text{ and } z' = f_{BI}(z'_{t+1})\}$ , that is,  $\mathcal{B}_{i,t}(z)$  is the set of terminal histories that (i) coincide with  $z$  up to date  $t$  and (ii) are reached by following the backward-induction choices from date  $t + 1$ . For example, for the game of Figure 2 (whose backward-induction solution is  $(a_1, a_2, (a_3, d_2))$  with corresponding terminal history  $a_1 a_2 a_3$ ) the model just described is shown in Figure 5.

By construction of the belief relations and by definition of backward-induction solution, at any state  $z$  and date  $t$ , if player  $i$  is active at  $z_t$  then he is rational there if and only if the action he takes there is the one prescribed by the

backward-induction solution, that is,  $z \in \mathbf{R}_t$  if and only if  $z_{t+1} = (f_{BI}(z_t))_{t+1}$ .<sup>14</sup> Let  $\hat{z}$  be the terminal history reached by the backward-induction solution, that is,  $\hat{z} = f_{BI}(\emptyset)$ . Then we have that  $\hat{z} \in B_t(\hat{z})$  for every date  $t \in T$  such that  $B_t(\hat{z}) \neq \emptyset$  and thus  $\hat{z} \in \mathbf{T}_t$  (that is, for every date  $t$ , the beliefs of the active player at  $\hat{z}$  are locally correct). Thus  $\hat{z} \in (\mathbf{T}_0 \cap \mathbf{R}_0) \cap B_0(\mathbf{T}_1 \cap \mathbf{R}_1) \cap B_0B_1(\mathbf{T}_2 \cap \mathbf{R}_2) \cap \dots \cap B_0B_1\dots B_{m-2}(\mathbf{T}_{m-1} \cap \mathbf{R}_{m-1})$ . ■

## References

- [1] R. Aumann. Backward induction and common knowledge of rationality. *Games and Economic Behavior*, 8:6–19, 1995.
- [2] R. Aumann. On the centipede game. *Games and Economic Behavior*, 23:97–105, 1998.
- [3] D. Balkenborg and E. Winter. A necessary and sufficient epistemic condition for playing backward induction. *Journal of Mathematical economics*, 27:325–345, 1997.
- [4] A. Baltag, S. Smets, and J. Zvesper. Keep ‘hoping’ for rationality: a solution to the backward induction paradox. *Synthese*, 169:301–333, 2009.
- [5] P. Battigalli and G. Bonanno. Recent results on belief, knowledge and the epistemic foundations of game theory. *Research in Economics*, 53:149–225, 1999.
- [6] P. Battigalli, A. Di-Tillio, and D. Samet. Strategies and interactive beliefs in dynamic games. Technical Report IGIER WP 375, Bocconi University, Milano, Italy, January 2011.
- [7] P. Battigalli and M. Siniscalchi. Strong belief and forward induction reasoning. *Journal of Economic Theory*, 106:356–391, 2002.
- [8] E. Ben-Porath. Nash equilibrium and backwards induction in perfect information games. *Review of Economic Studies*, 64:23–46, 1997.
- [9] A. Brandenburger. The power of paradox: some recent developments in interactive epistemology. *International Journal of Game Theory*, 35:465–492, 2007.
- [10] J. Halpern. Substantive rationality and backward induction. *Games and Economic Behavior*, 37:425–435, 2001.

---

<sup>14</sup>In the model of Figure 5  $\mathbf{R}_0 = \{a_1a_2b_3, a_1b_2, a_1a_2a_3\}$ ,  $\mathbf{R}_1 = \{a_1a_2b_3, a_1a_2a_3, b_1d_2\}$  and  $\mathbf{R}_2 = \{a_1a_2a_3, a_1b_2, b_1d_2, b_1c_2, b_1e_2\}$ . Thus  $B_0\mathbf{R}_1 = B_1\mathbf{R}_2 = B_0B_1\mathbf{R}_2 = Z$ . Note that  $a_1b_2 \in B_0\mathbf{R}_1$  but  $a_1b_2 \notin \mathbf{R}_1$  and thus at  $a_1b_2$  Player 1 has erroneous beliefs at date 0 about the rationality of the active player at date 1. In this model  $\mathbf{T}_0 = \{a_1a_2a_3, b_1d_2\}$ ,  $\mathbf{T}_1 = Z \setminus \{a_1a_2b_3\}$  and  $\mathbf{T}_2 = Z$ .

- [11] A. Perea. Epistemic foundations for backward induction: an overview. In J. van Benthem, D. Gabbay, and B. Löwe, editors, *Interactive logic. Proceedings of the 7th Augustus de Morgan Workshop*, volume 1 of *Texts in Logic and Games*, pages 159–193. Amsterdam University Press, 2007.
- [12] A. Perea. *Epistemic game theory: reasoning and choice*. Cambridge University Press, Cambridge, 2012.
- [13] D. Samet. Hypothetical knowledge and games with perfect information. *Games and Economic Behavior*, 17:230–251, 1996.
- [14] R. Stalnaker. Knowledge, belief and counterfactual reasoning in games. *Economics and Philosophy*, 12: 133–163, 1996.
- [15] R. Stalnaker. Belief revision in games: forward and backward induction. *Mathematical Social Sciences*, 36:31–56, 1998.
- [16] R. Stalnaker. On logics of knowledge and belief. *Philosophical Studies*, 128:169–199, 2006.