

Game Theoretic Decidability and Undecidability*

Tai-Wei Hu[†] and Mamoru Kaneko[‡]

05 September 2014, preliminary

Abstract

Logical inference is an engine for human thinking, especially, for decision making in an interdependent situation with more than one persons. We study the possibility of prediction/decision making in a finite 2-person game with pure strategies, following the Nash (-Johansen) (noncooperative solution) theory. Since some infinite regress naturally arises in this theory, we adopt a fixed-point extension EIR^2 of the epistemic logic KD^2 . The base logic KD^2 is adopted to capture individual decision making from the viewpoint of logical inference. Our results differ between solvable and unsolvable games. For the former, we have game theoretic decidability, i.e., player i can decide whether each of his strategies is a final decision or not. For the latter, he can neither decide it to be a possible decision nor can disprove it. This takes the form of Gödel's incompleteness theorem, while ours is a much simpler propositional theory. Our undecidability is related to "self-referential" as is Gödel's, but its main source is a discord generated by interdependence of payoffs and independent prediction/decision making.

Key words: Prediction/decision making, Infinite regress, Formal decidability, Undecidability, Incompleteness, Nash solution, Subsolution

1 Introduction

Logical inference is an engine for decision making in games with multiple players. Although game theory has studied decision making extensively, logical inference is kept informal. To study such a decision making process, we adopt a formal system of epistemic logic, the *epistemic infinite regress logic* EIR^2 . It is a fixed-point extension of the (propositional) epistemic logic KD^2 . Because of interdependence of players, prediction making is also required, and our logic allows us to model prediction making based on logical inference. At the same time, our approach emphasizes players' independence in terms of subjective thinking and this emphasis guides our choice of EIR^2 . The approach is coherent with Nash [16] and Johansen [9], who gave the noncooperative theory of prediction/decision making in a non-formalized manner. We study this theory in the logic EIR^2 .

We prove the game theoretic decidability and undecidability results, depending upon whether a game has the interchangeable set of Nash equilibria. The decidability result states that a player

*The authors are partially supported by Grant-in-Aids for Scientific Research No.26234567 and No.2312002, Ministry of Education, Science and Culture.

[†]Northwestern University, Illinois, USA, t-hu@kellogg.northwestern.edu

[‡]Faculty of Political Science and Economics, Waseda University, Tokyo, Japan, kaneko@sk.tsukuba.ac.jp

can reach a positive or a negative decision for each strategy, while the undecidability result states that for some strategy, he cannot reach either a positive or a negative decision.

Our approach takes various different perspectives from the standard literature of game theory as well as that of epistemic logic. Here, we explain those perspectives. For simplicity, we only focus on 2-person games, and a logic system with two players.

Fixed-point extension of KD^2 : The prediction/decision making process naturally leads to an infinite regress of beliefs. This regress begins subjectively in the mind of player i in his prediction making process, in which he simulates the other player's mind. The whole infinite regress arises recursively in a nested manner. In order to distinguish among the scopes of those minds, we adopt the epistemic logic KD^2 as the base logic. Adding the infinite regress operators, our fixed-point logic EIR^2 captures this infinite regress of beliefs¹.

As the concept of infinite regress of beliefs is closely related to the common knowledge, the logic EIR^2 is also related to the common knowledge logic CKL (cf., Fagin, *et al.* [4], and Meyer-van der Hoek [14]). In fact, if we add Axiom T(truthfulness) to the EIR^2 , then infinite regress collapse to common knowledge, and the resulting logic $EIR^2(T)$ becomes equivalent to CKL. Without Axiom T, EIR^2 can capture mutual subjectivity, which is not allowed in CKL.

Although some results in this paper are sharper in $EIR^2(T)$ than in EIR^2 , we take the latter as the basic system because $EIR^2(T)$ cannot capture the subjective nature of our problem but EIR^2 can.

Proof theory and model theory: Because of our focus on the prediction/decision making process with logical inference, we use a proof-theoretic system. We also use model theory (here, Kripke semantics) as a technical support, which is connected to our formal system via Kripke-soundness/completeness (see Hu-Kaneko [7]). We formalize a player's reasoning process in a formal system, instead of describing his mental states in a single (semantic) model^{2,3}. By soundness/completeness for EIR^2 , we can use Kripke models to evaluate provability via validity or finding a countermodel. In particular, the soundness part will be used to prove our game theoretic undecidability theorem.

Basic beliefs as non-logical axioms: As the formal Peano arithmetic is formulated by proper axioms in first-order classical logic, we postulate some basic beliefs as axioms for a player's prediction/decision making in the logic EIR^2 . Those basic beliefs include his understanding of the game and prediction/decision criterion. The inference from his beliefs to a decision is expressed as

$$\mathbf{B}_i(\Gamma_i^o) \vdash \mathbf{B}_i(I_i(s_i)). \quad (1)$$

That is, player i has basic beliefs Γ_i^o in his mind, and derives $I_i(s_i)$; his beliefs recommend s_i as a possible decision. The negative decision is described by $\mathbf{B}_i(\Gamma_i^o) \vdash \mathbf{B}_i(\neg I_i(s_i))$; his beliefs recommend him not to take s_i . Although (1) is expressed from the analyst's viewpoint, we intend to model these derivations as occurring in player i 's mind. In fact, in the logic EIR^2 , $\mathbf{B}_i(\Gamma_i^o) \vdash$

¹Alternatively, we can adopt an infinitary logic. Hu *et al.* [8].discusses relationships between the IER^2 and its infinitary counterpart.

²The model-theoretic standpoint has been taken almost exclusively in the literature of epistemic logic with applications to game theory; for example, see van Benthem *et al.* [20], the various papers in Brandenburger [3], and van Benthem [19]. Some exceptions are Kaneko-Nagashima [10], Kline [13], and Suzuki [18], where the proof-theoretic standpoint is taken.

³Many aspects involved in playing a game are considered in van Benthem *et al.* [20] and van Benthem [19]. In particular, matrix games are formulated by means of logic in Chap.12 of [19]. Nevertheless, an individual thought process of prediction/decision making is only indirectly treated.

$\mathbf{B}_i(\mathbf{I}_i(s_i))$ ($\mathbf{B}_i(\Gamma_i^o) \vdash \mathbf{B}_i(\neg\mathbf{I}_i(s_i))$) is equivalent to $\Gamma_i^o \vdash \mathbf{I}_i(s_i)$ ($\Gamma_i^o \vdash \neg\mathbf{I}_i(s_i)$); this equivalence is formally stated in Lemma 2.5, and we interpret there the latter as occurring in player i 's mind. The choice of the base logic KD^2 is essential for this equivalence.

Game theoretic concepts: We only consider finite 2-person strategic games with pure strategies. This simple setting is rich enough to obtain both decidability and undecidability results. In fact, the characterization of games with decidability/undecidability corresponds to the interchangeability requirement in Nash [16]. Interchangeability captures players' independence in *ex ante* prediction/decision making, but Nash did not make a formal distinction between prediction and decision. Johansen [9] discussed Nash's theory in a more philosophical manner with a conceptual distinction between prediction and decision. As our axioms for prediction/decision making formalize his argument in the logic EIR^2 , the resulting system is called the *formalized Nash-Johansen theory* (for short, *the formalized Nash theory*).

Axiomatic formulation of prediction/decision making: We postulate three axioms, N0_i , N1_i , and N2_i , given in Section 4, for prediction/decision making. They are in the scope of the mind of player i , expressed as $\mathbf{B}_i(\text{N012}_i) := \mathbf{B}_i(\text{N0}_i \wedge \text{N1}_i \wedge \text{N2}_i)$. To make his prediction about player j 's decision, player i uses the belief $\mathbf{B}_i\mathbf{B}_j(\text{N012}_j)$, where N012_j is the same as N012_i with the replacement of i with j . For the same reason, $\mathbf{B}_i\mathbf{B}_j(\text{N012}_j)$ requires $\mathbf{B}_i\mathbf{B}_j\mathbf{B}_i(\text{N012}_i)$, and so on. Therefore, to complete prediction making, player i would meet an infinite regress of beliefs:

$$\mathbf{B}_i(\text{N012}_i), \mathbf{B}_i\mathbf{B}_j(\text{N012}_j), \mathbf{B}_i\mathbf{B}_j\mathbf{B}_i(\text{N012}_i), \dots \quad (2)$$

This is captured by the fixed-point operator, $\mathbf{I}_i(\text{N012}_i; \text{N012}_j)$, in the logic EIR^2 .

The infinite sequence (2), *a fortiori*, $\mathbf{I}_i(\text{N012}_i; \text{N012}_j)$, has a self-referential structure: Itself occurs in the scope of $\mathbf{B}_i(\cdot)$, the counterpart for player j is in the scope of $\mathbf{B}_i\mathbf{B}_j(\cdot)$, and itself occurs again in $\mathbf{B}_i\mathbf{B}_j\mathbf{B}_i(\cdot)$, and so on. This self-referential structure is crucial for our undecidability result.

Conceptually, the infinite regress, $\mathbf{I}_i(\text{N012}_i; \text{N012}_j)$, is our basic postulate for prediction/decision making. Mathematically, however, it only provides a necessary condition for possible decisions. We formulate another axiom (schema), $\mathbf{I}_i(\mathbf{WF})$, that gives the sufficiency of this postulate to determine a possible decision.

Formalized Nash theory: The set of beliefs $\mathbf{I}_i(\text{N012}_i; \text{N012}_j), \mathbf{I}_i(\mathbf{WF})$ describes prediction/decision making without concrete information about the game situation. We formulate the basic beliefs of the game situation (including strategies and payoffs) by $\mathbf{I}_i(\mathbf{g}) := \mathbf{I}_i(g_i; g_j)$. This addition completes our postulates of player i 's basic beliefs: $\Delta_i(\mathbf{g}) = \{\mathbf{I}_i(\mathbf{g}), \mathbf{I}_i(\text{N012}_i; \text{N012}_j)\} \cup \mathbf{I}_i(\mathbf{WF})$, which plays the role of $\mathbf{B}_i(\Gamma_i^o)$ in (1). Note that the set of beliefs $\Delta_i(\mathbf{g})$ depends upon the game situation $\mathbf{g} = (g_i; g_j)$. The pair $(\text{EIR}^2; \Delta_i(\mathbf{g}))$ forms the formalized Nash theory for \mathbf{g} .

The literature of game theory tends to focus on the resulting outcome(s) from a solution/equilibrium theory. In our context, this can be stated as the following question:

(i): What decisions and predictions does $(\text{EIR}^2; \Delta_i(\mathbf{g}))$ recommend?

This question presumes that the theory $(\text{EIR}^2; \Delta_i(\mathbf{g}))$ has recommendations. However, we should ask the following question in the first place.

(ii): Does $(\text{EIR}^2; \Delta_i(\mathbf{g}))$ recommend any decision?

In fact, the answer to the second question is related to Nash's [16] interchangeability condition.

We say that a game is *solvable* when the set of Nash equilibria is interchangeable, i.e., the set has a product structure. Here, we give three examples of games; two are solvable and one is not. In Table 1.1, each player has three strategies, and his payoff is given in the matrix (the first and second entries are players 1's and 2's payoffs). The superscript NE stands for Nash equilibrium, explained in Section 3. Table 1.1 has a unique Nash equilibrium. Table 1.2, called the *battle of the sexes*, has two Nash equilibria; these are not solvable because the set is not a product set. Table 1.3, called the *matching pennies*, has the empty set of Nash equilibria. Tables 1.1 and 1.3 are solvable games.

	s ₂₁	s ₂₂	s ₂₃
s ₁₁	2, 4	2, 2	4, 0
s ₁₂	3, 3 ^{NE}	4, 2	3, 0
s ₁₃	0, 0	5, 5	2, 6

	s ₂₁	s ₂₂
s ₁₁	2, 1 ^{NE}	0, 0
s ₁₂	0, 0	1, 2 ^{NE}

	s ₂₁	s ₂₂
s ₁₁	1, -1	-1, 1
s ₁₂	-1, 1	1, -1

Positive, negative decisions, and undecidable: Our results answer the question (ii) as follows. When a game is solvable, we have the decidability result: for *any* strategy s_i for player i ,

$$\text{either } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i)) \text{ or } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg\mathbf{I}_i(s_i)). \quad (3)$$

For Table 1.1, the set of beliefs $\Delta_1(\mathbf{g})$ recommends to player 1 to take \mathbf{s}_{12} as a positive decision but not to take either \mathbf{s}_{11} or \mathbf{s}_{13} . In Table 1.3, $\Delta_1(\mathbf{g})$ recommends all as negative decisions.

The main result of the paper shows that when a game \mathbf{g} is not solvable such as Table 1.2, there is *some* strategy s_i for each player i such that

$$\text{neither } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i)) \text{ nor } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg\mathbf{I}_i(s_i)). \quad (4)$$

That is, player i cannot decide with the belief set $\Delta_i(\mathbf{g})$ whether s_i is a positive or negative decision. In Table 1.2, this holds for both strategies. This situation differs entirely from the case where $\Delta_i(\mathbf{g})$ gives negative recommendations for all strategies as in Table 1.3; in the latter case, he may look for a different way for decision making, but in the former, i.e., (4), he may not be able to notice this undecidability itself, and get stuck in his decision making.

Relations to Gödel's incompleteness theorem and the source for our undecidability: The result (4) has the same form as Gödel's incompleteness theorem (cf., Boolos [2], Mendelson [15]), but both interpretation and source for incompleteness differ. Gödel's theorem is about the Peano Arithmetic and based on the self-referential structure. Although the self-referential structure in the infinite regress of beliefs is crucial to our undecidability result, it is not the only source. Our answer to the question (ii) above reveals that the basic belief $\mathbf{I}_i(\mathbf{g})$ plays an indispensable role. Among the three components of $\Delta_i(\mathbf{g})$, $\mathbf{I}_i(\mathbf{N}012_i; \mathbf{N}012_j)$ and $\mathbf{I}_i(\mathbf{WF})$ are symmetric between the two players, but discords included in $\mathbf{I}_i(\mathbf{g})$ may bring about undecidability. A detailed comparison with Gödel's theorem will be given in Section 6.

The format of the paper is as follows: Section 2 formulates the logic EIR². Section 3 gives various game theoretic concepts. Section 4 gives three axioms for prediction/decision making, and the decidability result for a solvable game. Section 5 presents the undecidability result for an unsolvable game. Section 6 gives concluding remarks.

2 The Epistemic Infinite Regress Logic EIR²

We formulate the logic EIR² with the language for 2-person strategic games in Sections 2.1, 2.2, and give its semantics in Section 2.3. The language presumes the sets of strategies but this restriction is not essential for our argument.

2.1 Language

Let S_i be a nonempty finite *strategy* set for player $i = 1, 2$. We adopt the atomic formulae:

atomic preference formulae: $\text{Pr}_i(s; t)$ for $i = 1, 2$ and $s, t \in S = S_1 \times S_2$;

atomic decision formulae: $\text{I}_i(s_i)$ for $s_i \in S_i$, $i = 1, 2$.

The atomic formula $\text{Pr}_i(\cdot; \cdot)$ expresses the preference relation of player i ; $\text{Pr}_i(s; t)$ means that player i *weakly prefers* the strategy pair $s = (s_1, s_2)$ to the pair $t = (t_1, t_2)$. The atomic formula $\text{I}_i(s_i)$ expresses the idea that, from player i 's perspective, s_i is a *possible final decision* for him.

Now we proceed to have logical connectives and epistemic operators:

logical connective symbols: \neg (not), \supset (imply), \wedge (and), \vee (or);⁴

unary belief operators: $\mathbf{B}_1(\cdot)$, $\mathbf{B}_2(\cdot)$; *binary infinite-regress operators:* $\mathbf{Ir}_1(\cdot, \cdot)$, $\mathbf{Ir}_2(\cdot, \cdot)$;

parentheses: $(,)$.

We stipulate that j refers to the other player than i . Player i 's prediction about player j 's decision is expressed as $\mathbf{B}_j(\text{I}_j(s_j))$, but this should occur in the scope of $\mathbf{B}_i(\cdot)$. We use a pair of formulae, (A_1, A_2) , as arguments of the binary operators $\mathbf{Ir}_1(\cdot, \cdot)$ and $\mathbf{Ir}_2(\cdot, \cdot)$, and the intended meaning of the formula $\mathbf{Ir}_i(A_1, A_2)$ is that player i 's subjective belief of the infinite regress of beliefs about A_i and A_j . We write $\mathbf{Ir}_i(A_1, A_2)$ also as $\mathbf{Ir}_i(A_i; A_j)$ and sometimes $\mathbf{Ir}_i[A_i; A_j]$.

We define the sets of *formulae*, denoted by \mathcal{P} , by the following induction:

- (o) all atomic formulae are formulae;
- (i) if A, B are formulae, then so are $(A \supset B)$, $(\neg A)$, $\mathbf{B}_i(A)$ for $i = 1, 2$;
- (ii) if $\mathbf{A} = (A_1, A_2)$ is a pair of formulae, then $\mathbf{Ir}_i(\mathbf{A})$ is also a formula;
- (iii) if Φ is a finite (nonempty) set of formulae, then $(\wedge \Phi)$ and $(\vee \Phi)$ are formulae⁵.

We say that a formula A is *non-epistemic* iff $\mathbf{B}_i(\cdot)$ or $\mathbf{Ir}_i(\cdot, \cdot)$ does not occur in A for $i = 1, 2$. The set of nonepistemic formulae is denoted by \mathcal{P}_N . We say that A_i is a *game formula for i* iff it contains atomic formulae of the form $\text{Pr}_i(\cdot; \cdot)$ only, that is, no occurrences of $\text{Pr}_j(\cdot; \cdot)$, $\text{I}_i(\cdot)$, and $\text{I}_j(\cdot)$; and that A is a *game formula* iff the atomic formulae occurring in A are of the form $\text{Pr}_1(\cdot; \cdot)$ or $\text{Pr}_2(\cdot; \cdot)$. A game formula expresses a reality of the target situation together with, potentially, beliefs about them, while the atomic decision formulae $\text{I}_i(s_i)$'s are used to describe a player's thinking about prediction/decision making.

We write $\wedge\{A, B\}$, $\wedge\{A, B, C\}$ as $A \wedge B$, $A \wedge B \wedge C$, etc., and $(A \supset B) \wedge (B \supset A)$ as $A \equiv B$.

⁴Since we adopt classical logic as the base logic, we can abbreviate some of those connectives. Since, however, our aim is to study logical inference for decision making rather than semantic contents, we use a full system.

⁵We presume the identity of finite sets in our language.

We abbreviate parentheses or use different ones such as $[,]$ when no confusions are expected.

2.2 Proof theory of EIR²

We start with an explicit formulation of classical logic, which consists of five axiom (schemata) and three inference rules: for all formulae A, B, C , and finite nonempty sets Φ of formulae,

- L1** $A \supset (B \supset A)$;
- L2** $(A \supset (B \supset C)) \supset ((A \supset B) \supset (A \supset C))$;
- L3** $(\neg A \supset \neg B) \supset ((\neg A \supset B) \supset A)$;
- L4** $\wedge \Phi \supset A$, where $A \in \Phi$;
- L5** $A \supset \vee \Phi$, where $A \in \Phi$;

$$\frac{A \supset B \quad A}{B} \text{MP} \quad \frac{\{A \supset B : B \in \Phi\}}{A \supset \wedge \Phi} \wedge\text{-rule} \quad \frac{\{B \supset A : B \in \Phi\}}{\vee \Phi \supset A} \vee\text{-rule}.$$

The epistemic logic KD² is defined by adding, to classical logic, two epistemic axioms and one inference rule for the belief operators $\mathbf{B}_i(\cdot)$: for all formulae A, C , and for $i = 1, 2$,

- K** $\mathbf{B}_i(A \supset C) \supset (\mathbf{B}_i(A) \supset \mathbf{B}_i(C))$;
- D** $\neg \mathbf{B}_i(\neg A \wedge A)$;
- Necessitation** $\frac{A}{\mathbf{B}_i(A)}$.

Then, we have the *epistemic infinite regress logic* EIR², by adding one axiom (schema) and one inference rule for the infinite regress operators $\mathbf{I}_i(\cdot, \cdot)$: For $i = 1, 2$, and two pairs of formulae $\mathbf{A} = (A_1, A_2)$, $\mathbf{D} = (D_1, D_2)$,

- IRA_i** $\mathbf{I}_i(\mathbf{A}) \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j(A_j) \wedge \mathbf{B}_i \mathbf{B}_j(\mathbf{I}_i(\mathbf{A}))$;
- IRI_i** $\frac{D_i \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j(A_j) \wedge \mathbf{B}_i \mathbf{B}_j(D_i)}{D_i \supset \mathbf{I}_i(\mathbf{A})}$.

Axiom IRA_i has a fixed-point structure in the sense that $\mathbf{B}_i \mathbf{B}_j(\mathbf{I}_i(\mathbf{A}))$ appears as an implication of $\mathbf{I}_i(\mathbf{A})$. Replacing $\mathbf{I}_i(\mathbf{A})$ in $\mathbf{B}_i \mathbf{B}_j(\mathbf{I}_i(\mathbf{A}))$ with its implication $\mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j(A_j)$ (formally with K and Nec), $\mathbf{I}_i(\mathbf{A})$ implies the following infinite regress of beliefs:

$$\{\mathbf{B}_i(A_i), \mathbf{B}_i \mathbf{B}_j(A_j), \mathbf{B}_i \mathbf{B}_j \mathbf{B}_i(A_i), \dots\}. \quad (5)$$

Rule IRI_i states that $\mathbf{I}_i(\mathbf{A})$ is the logically weakest formula satisfying the property described in IRA_i, that is, if D_i enjoys it, then D_i implies $\mathbf{I}_i(\mathbf{A})$. Our completeness-soundness (Theorem 2.1) shows that $\mathbf{I}_i(\mathbf{A})$ captures faithfully the set in (5).

A *proof* $P = \langle X, <; \psi \rangle$ consists of a finite tree $\langle X, < \rangle$ and a function $\psi : X \rightarrow \mathcal{P}$ with the following requirements:

- P1** for each node $x \in X$, $\psi(x)$ is a formula attached to x ;
- P2** for each leaf x in $\langle X, < \rangle$, $\psi(x)$ is an instance of the axiom schemata;

P3 for each non-leaf x in $\langle X, < \rangle$,

$$\frac{\{\psi(y) : y \text{ is an immediate predecessor of } x\}}{\psi(x)}$$

is an instance of the above five inference rules.

We call P a *proof of* A iff $\psi(x_0) = A$, where x_0 is the root of $\langle X, < \rangle$. We say that A is *provable*, denoted by $\vdash A$, iff there is a proof of A . For a set of formulae Γ , we write $\Gamma \vdash A$ iff $\vdash A$ or there is a finite nonempty subset Φ of Γ such that $\vdash \wedge \Phi \supset A$. This treatment of non-logical assumptions is crucial in our study⁶.

The following are basic to classical logic and/or KD^2 (cf., Kaneko [11]). We use them without referring.

Lemma 2.1. *Let $A \in \mathcal{P}$, Φ a finite set of formulae, and $i = 1, 2$. Then, (1) $\vdash A \supset B$ and $\vdash B \supset C$ imply $\vdash A \supset C$; (2) $\vdash (A \wedge B \supset C) \equiv (A \supset (B \supset C))$; (3) $\vdash \mathbf{B}_i(\neg A) \supset \neg \mathbf{B}_i(A)$; (4) $\vdash \vee \mathbf{B}_i(\Phi) \supset \mathbf{B}_i(\vee \Phi)$; (5) $\vdash \mathbf{B}_i(\wedge \Phi) \equiv \wedge \mathbf{B}_i(\Phi)$.*

We will use the following three lemmas in the subsequent discussions. First, from Axiom IRA_i and Rule IRI_i ($i = 1, 2$), the operators $\mathbf{I}r_i(\cdot, \cdot)$ and $\mathbf{I}r_j(\cdot, \cdot)$ may appear to be independent of one another, but they are interdependent.

Lemma 2.2. (Epistemic content) *Let $\mathbf{A} = (A_1, A_2)$ be a pair of formulae. Then, $\vdash \mathbf{I}r_i(\mathbf{A}) \equiv \mathbf{B}_i(A_i \wedge \mathbf{I}r_j(\mathbf{A}))$ for $i = 1, 2$.*

Proof. Let us see $\vdash \mathbf{B}_i(A_i \wedge \mathbf{I}r_j(\mathbf{A})) \supset \mathbf{I}r_i(\mathbf{A})$. Let $D_i = \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A}))$ for $i = 1, 2$. By IRA_j (and, Nec, K), we have $\vdash D_i \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j(A_j) \wedge \mathbf{B}_i \mathbf{B}_j \mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A}))$. Since the last two conjuncts are equivalent to $\mathbf{B}_i \mathbf{B}_j(D_i)$, we have $\vdash D_i \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i \mathbf{B}_j(A_j) \wedge \mathbf{B}_i \mathbf{B}_j(D_i)$. Using IRI_i , we have $\vdash \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A})) \supset \mathbf{I}r_i(\mathbf{A})$.

The above for j implies $\vdash \mathbf{B}_i(D_j) \supset \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A}))$. Hence, $\vdash \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(D_j) \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A}))$. Since $\vdash \mathbf{I}r_i(\mathbf{A}) \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(D_j)$ by IRA_i , we have $\vdash \mathbf{I}r_i(\mathbf{A}) \supset \mathbf{B}_i(A_i) \wedge \mathbf{B}_i(\mathbf{I}r_j(\mathbf{A}))$. ■

This lemma enables us to talk about the *epistemic content* of $\mathbf{I}r_i(\mathbf{A})$;

$$\mathbf{I}r_i^o(\mathbf{A}) := A_i \wedge \mathbf{I}r_j(\mathbf{A}), \tag{6}$$

which plays a crucial role in our consideration of prediction/decision making.

Lemma 2.3. (Basic properties for $\mathbf{I}r_i(\cdot, \cdot)$) *Let $\mathbf{A} = (A_1, A_2)$ and $\mathbf{C} = (C_1, C_2)$ be two pairs of formulae in \mathcal{P} and $i = 1, 2$.*

(1) *If $\vdash \mathbf{I}r_k(\mathbf{A}) \supset \mathbf{B}_k(C_k)$ for $k = 1, 2$, then $\vdash \mathbf{I}r_i(\mathbf{A}) \supset \mathbf{I}r_i(\mathbf{C})$. In particular, if $\vdash C_k$ for $k = 1, 2$, then $\vdash \mathbf{I}r_i(\mathbf{C})$.*

(2) $\vdash \mathbf{I}r_i(\mathbf{A}) \supset \mathbf{I}r_i(\mathbf{I}r_1^o(\mathbf{A}), \mathbf{I}r_2^o(\mathbf{A}))$;

(3) $\vdash \mathbf{I}r_i(A_1 \wedge C_1, A_2 \wedge C_2) \equiv \mathbf{I}r_i(\mathbf{A}) \wedge \mathbf{I}r_i(\mathbf{C})$;

(4) $\vdash \mathbf{I}r_i(A_1 \supset C_1, A_2 \supset C_2) \supset (\mathbf{I}r_i(\mathbf{A}) \supset \mathbf{I}r_i(\mathbf{C}))$;

(5) $\vdash \mathbf{I}r_i(\neg A_i; A_j) \supset \neg \mathbf{I}r_i(\mathbf{A})$, $\vdash \mathbf{I}r_i(A_i; \neg A_j) \supset \neg \mathbf{I}r_i(\mathbf{A})$, and $\vdash \mathbf{I}r_i(\neg A_i; \neg A_j) \supset \neg \mathbf{I}r_i(\mathbf{A})$.

⁶Since the deduction theorem (cf., Mendelson [15]) does not hold in epistemic logic, the introduction of non-logical axioms differs from in classical logic. We adopt the classical manner.

Proof. (1): Let $\vdash \mathbf{Ir}_k(\mathbf{A}) \supset \mathbf{B}_k(C_k)$ for $k = 1, 2$. We show $\vdash \mathbf{Ir}_i(\mathbf{A}) \supset \mathbf{B}_i(C_i) \wedge \mathbf{B}_i\mathbf{B}_j(C_j) \wedge \mathbf{B}_i\mathbf{B}_j(\mathbf{Ir}_i(\mathbf{A}))$, which implies, by IRI_i , $\vdash \mathbf{Ir}_i(\mathbf{A}) \supset \mathbf{Ir}_i(\mathbf{C})$. First, $\vdash \mathbf{B}_i(\mathbf{Ir}_j(\mathbf{A})) \supset \mathbf{B}_i\mathbf{B}_j(C_j)$ by Nec and K. By Lemma 2.2, we have $\vdash \mathbf{Ir}_i(\mathbf{A}) \supset \mathbf{B}_i\mathbf{B}_j(C_j)$. By IRA_i , we have $\vdash \mathbf{Ir}_i(\mathbf{A}) \supset \mathbf{B}_i\mathbf{B}_j(\mathbf{Ir}_i(\mathbf{A}))$. By \wedge -rule, we have the target.

The other claims (2)-(4) follow (1). Here, we show (3). Since $\vdash \mathbf{Ir}_k(A_1 \wedge C_1, A_2 \wedge C_2) \supset \mathbf{B}_k(A_k)$ for $k = 1, 2$, we have, by (1), $\vdash \mathbf{Ir}_k(A_1 \wedge C_1, A_2 \wedge C_2) \supset \mathbf{Ir}_i(\mathbf{A})$. Similarly, $\vdash \mathbf{Ir}_k(A_1 \wedge C_1, A_2 \wedge C_2) \supset \mathbf{Ir}_i(\mathbf{C})$. Hence, we have the one direction. Consider the converse. We have $\vdash \mathbf{Ir}_k(\mathbf{A}) \wedge \mathbf{Ir}_k(\mathbf{C}) \supset \mathbf{B}_k(A_k \wedge C_k)$ for $k = 1, 2$. We have $\vdash \mathbf{Ir}_i(\mathbf{A}) \wedge \mathbf{Ir}_i(\mathbf{C}) \supset \mathbf{B}_i\mathbf{B}_j(A_j \wedge C_j)$, and $\vdash \mathbf{Ir}_i(\mathbf{A}) \wedge \mathbf{Ir}_i(\mathbf{C}) \supset \mathbf{B}_i\mathbf{B}_j(\mathbf{Ir}_i(\mathbf{A}) \wedge \mathbf{Ir}_i(\mathbf{C}))$. Then, by IRI_i , $\vdash \mathbf{Ir}_i(\mathbf{A}) \wedge \mathbf{Ir}_i(\mathbf{C}) \supset \mathbf{Ir}_i(A_1 \wedge C_1, A_2 \wedge C_2)$.

(5): Consider only the first one. Since $\vdash \mathbf{Ir}_i(\neg A_i; A_j) \supset \mathbf{B}(\neg A_i)$, we have $\vdash \mathbf{Ir}_i(\neg A_i; A_j) \supset \neg \mathbf{B}(A_i)$. Then, using the contrapositive of IRA_i , i.e., $\vdash \neg[\mathbf{B}_i(A_i) \wedge \mathbf{B}_i\mathbf{B}_j(A_j) \wedge \mathbf{B}_i\mathbf{B}_j(\mathbf{Ir}_i(\mathbf{A}))] \supset \neg \mathbf{Ir}_i(\mathbf{A})$, we have $\vdash \mathbf{Ir}_i(\neg A_i; A_j) \supset \neg \mathbf{Ir}_i(\mathbf{A})$. ■

The following statements for $\mathbf{Ir}_i^o(\cdot; \cdot)$ correspond to IRA_i and IRI_i for $\mathbf{Ir}_i(\cdot; \cdot)$.

Lemma 2.4. (Admissible formulae and inference) Let $\mathbf{A} = (A_i; A_j)$ and D_i be any formulae. Then,

(1) $(\text{IRA}_i^o) \vdash \mathbf{Ir}_i^o(\mathbf{A}) \supset A_i \wedge \mathbf{B}_j(A_j) \wedge \mathbf{B}_j\mathbf{B}_i(\mathbf{Ir}_i^o(\mathbf{A}))$;

(2) (IRI_i^o) If $\vdash D_i \supset A_i \wedge \mathbf{B}_j(A_j) \wedge \mathbf{B}_j\mathbf{B}_i(D_i)$, then $\vdash D_i \supset \mathbf{Ir}_i^o(A_i; A_j)$.

Proof. (1): By (6), $\vdash \mathbf{Ir}_i^o(\mathbf{A}) \supset A_i \wedge \mathbf{Ir}_j(\mathbf{A})$. By Lemma 2.2 for j , we have $\vdash \mathbf{Ir}_i^o(\mathbf{A}) \supset A_i \wedge \mathbf{B}_j(A_j) \wedge \mathbf{B}_j(\mathbf{Ir}_i(\mathbf{A}))$.

(2): Let $\vdash D_i \supset A_i \wedge \mathbf{B}_j(A_j) \wedge \mathbf{B}_j\mathbf{B}_i(D_i)$. Since $\vdash D_i \supset \mathbf{B}_j\mathbf{B}_i(D_i)$ and $\vdash D_i \supset A_i$, we have $\vdash D_i \supset \mathbf{B}_j\mathbf{B}_i(A_i)$. Thus, $\vdash D_i \supset \mathbf{B}_j(A_j) \wedge \mathbf{B}_j\mathbf{B}_i(A_i) \wedge \mathbf{B}_j\mathbf{B}_i(D_i)$. By IRI_i , we have $\vdash D_i \supset \mathbf{Ir}_j(A_i; A_j)$. Thus, $\vdash D_i \supset A_i \wedge \mathbf{Ir}_j(A_i; A_j)$, which is $\vdash D_i \supset \mathbf{Ir}_i^o(A_i; A_j)$ by (6). ■

The main undecidability result of the paper holds in stronger systems than EIR^2 , such as those obtained from EIR^2 by adding Axiom T (truthfulness): $\mathbf{B}_i(A) \supset A$; Axiom 4 (positive introspection): $\mathbf{B}_i(A) \supset \mathbf{B}_i\mathbf{B}_i(A)$; and/or Axiom 5 (negative introspection): $\neg \mathbf{B}_i(A) \supset \mathbf{B}_i(\neg \mathbf{B}_i(A))$. We choose KD^2 as the base logic to give a clear-cut description of each player's logical inference. This is stated by the scope lemma (Lemma 2.5), which would not hold in any stronger system mentioned above⁷.

Nevertheless, Axiom T helps us understand the fixed-point formula $\mathbf{Ir}_i(\mathbf{A})$. Now, let us see the common knowledge logic CKL (cf., Fagin *et al.* [4] and Meyer-van der Hoek [14]). The logic CKL uses only one operator, $\mathbf{C}(\cdot)$, and adds the following axiom and rule to KD^2 :

CKA: $\mathbf{C}(A) \supset A \wedge \mathbf{B}_1(\mathbf{C}(A)) \wedge \mathbf{B}_2(\mathbf{C}(A))$;

CKI: $\frac{D \supset A \wedge \mathbf{B}_1(D) \wedge \mathbf{B}_2(D)}{D \supset \mathbf{C}(A)}$.

Axiom CKA and Rule CKI are interpreted as meaning that $\mathbf{C}(A)$ describes the common knowledge of A from the outside analyst's perspective. In contrast, $\mathbf{Ir}_i(\mathbf{A})$ describes player i 's beliefs from his subjective perspective. This difference is reflected by the counterpart of (5) in CKL,

⁷We regard KD^2 as the basic system; Axiom K and Necessitation give the inference ability of classical logic to each player. If Axiom D is dropped, player's beliefs can be arbitrary with no restrictions; for instance, $\mathbf{B}_i(p) \not\vdash \neg \mathbf{B}_i(\neg p)$ holds. Axiom D avoids this contradictory beliefs.

i.e., $\mathbf{C}(A)$ captures the entire set:

$$\{A, \mathbf{B}_1(A), \mathbf{B}_2(A), \mathbf{B}_1\mathbf{B}_2(A), \mathbf{B}_2\mathbf{B}_1(A), \mathbf{B}_1\mathbf{B}_2\mathbf{B}_2(A), \dots\}. \quad (7)$$

This set of formulae having all finite sequences of $\mathbf{B}_2\mathbf{B}_1\dots$ including the repetitive ones such as $\mathbf{B}_1\mathbf{B}_2\mathbf{B}_2$, while each in (5) has the outer $\mathbf{B}_i(\cdot)$ and all $\mathbf{B}_i\mathbf{B}_j\dots$ are alternating.

If we add Axiom T to the logic EIR^2 , which is denoted by $\text{EIR}^2(\text{T})$, an infinite regress collapses to common knowledge. Lemma 2.2 implies $\vdash \mathbf{I}r_i(A_1, A_2) \equiv \mathbf{I}r_j(A_1, A_2) (\equiv \mathbf{I}r_i^o(A_1, A_2))$ for $i = 1, 2$ in $\text{EIR}^2(\text{T})$, and for any formulae A_1, A_2 and D ,

$$\mathbf{cka}: \vdash \mathbf{I}r_i(A_1, A_2) \supset (A_1 \wedge A_2) \wedge \mathbf{B}_1\mathbf{I}r_i(A_1, A_2) \wedge \mathbf{B}_2\mathbf{I}r_i(A_1, A_2);$$

$$\mathbf{cki}: \text{if } \vdash D \supset (A_1 \wedge A_2) \wedge \mathbf{B}_1(D) \wedge \mathbf{B}_2(D), \text{ then } \vdash D \supset \mathbf{I}r_i(A_1, A_2).$$

Thus, in $\text{EIR}^2(\text{T})$, CKA and CKI are derived formulae and admissible rule for $\mathbf{I}r_i(A_1, A_2)$, that is, $\mathbf{I}r_i(A_1, A_2)$ means the common knowledge of $A_1 \wedge A_2$.

We will use the *belief eraser* ε_0 : the nonepistemic formula $\varepsilon_0(A) \in \mathcal{P}_N$ is obtained from $A \in \mathcal{P}$ by eliminating all occurrences of $\mathbf{B}_1(\cdot), \mathbf{B}_2(\cdot)$ and replacing $\mathbf{I}r_i(A_1, A_2)$ by $\varepsilon_0(A_1) \wedge \varepsilon_0(A_2)$. Then, we have

$$\vdash A \text{ implies } \vdash_0 \varepsilon_0(A), \quad (8)$$

where \vdash_0 is the provability relation of classical logic in \mathcal{P}_N . This is proved by induction on a proof of A from its leaves (cf., Kaneko-Nagashima [10]).

2.3 Kripke semantics and the soundness/completeness of EIR^2

Here, we report soundness/completeness for EIR^2 with respect to the Kripke semantics. We use the soundness part for the main undecidability result.

A Kripke frame $\langle W; R_1, R_2 \rangle$ consists of a nonempty set W of possible worlds and an accessibility relation R_i for player $i = 1, 2$. We say that a frame $\langle W; R_1, R_2 \rangle$ is *serial* iff for $i = 1, 2$ and for all $w \in W$, $wR_i u$ for some $u \in W$. A *truth assignment* τ is a function from $W \times AF$ to $\{\top, \perp\}$, where AF is the set of atomic formulae. A pair $M = (\langle W; R_1, R_2 \rangle, \tau)$ is called a *model*. When $\langle W; R_1, R_2 \rangle$ is serial, we say that M is a serial model.

We say that $\langle (w_0, i_0), \dots, (w_\nu, i_\nu), w_{\nu+1} \rangle$ ($\nu \geq 0$) is an *alternating chain* iff $i_{k-1} \neq i_k$ for $k = 1, \dots, \nu$ and $w_{k-1}R_{i_{k-1}}w_k$ for $k = 1, \dots, \nu + 1$. The alternating structure corresponds to the set given by (5). This is used for evaluating the truth values of formulae $\mathbf{I}r_i(A_1, A_2)$, $i = 1, 2$.

The valuation in (M, w) , denoted by $(M, w) \models$, is defined over \mathcal{P} by induction on the length of a formula as follows:

$$\mathbf{V0} \text{ for any } A \in AF, (M, w) \models A \iff \tau(w, A) = \top;$$

$$\mathbf{V1} (M, w) \models \neg A \iff (M, w) \not\models A;$$

$$\mathbf{V2} (M, w) \models A \supset B \iff (M, w) \not\models A \text{ or } (M, w) \models B;$$

$$\mathbf{V3} (M, w) \models \wedge \Phi \iff (M, w) \models A \text{ for all } A \in \Phi;$$

$$\mathbf{V4} (M, w) \models \vee \Phi \iff (M, w) \models A \text{ for some } A \in \Phi;$$

$$\mathbf{V5} (M, w) \models \mathbf{B}_i(A) \iff (M, v) \models A \text{ for all } v \text{ with } wR_iv;$$

V6 $(M, w) \models \mathbf{Ir}_i(A_1, A_2) \iff (M, w_{\nu+1}) \models A_{i_\nu}$ for any alternating chain $\langle (w_0, i_0), \dots, (w_\nu, i_\nu), w_{\nu+1} \rangle$ with $(w_0, i_0) = (w, i)$.

The steps other than V6 are standard. V6 is similar to the valuation for the common knowledge operator in CKL; the only difference is to use alternating reachability for two formulae, instead of simple reachability (cf., Fagin *et al.* [4], Meyer-van der Hoek [14]).

We have the following soundness/completeness theorem.

Theorem 2.1. (Soundness and Completeness) *Let $A \in \mathcal{P}$. Then, $\vdash A$ in EIR^2 if and only if $(M, w) \models A$ for all serial models $M = (\langle W; R_1, R_2 \rangle, \tau)$ and any $w \in W$.*

Soundness (only-if) will be used to prove our undecidability result (Theorem 5.1). It is proved as follows: Let $P = (X, <; \psi)$ be a proof of A . Then, by induction on the tree structure of $(X, <)$ from its leaves, we show that for any $x \in X$, $\vdash \psi(x)$ implies $\models \psi(x)$. The two new steps are : (1) $\models C$ for any instance C of IRA_i ; and (2) the validity relation \models preserves Rule IRI_i . Both steps follow from V6. The proof of completeness is given in Hu-Kaneko [7], which also shows that the theorem still holds under any additions of Axioms T, 4 and 5.

Theorem 2.1 shows that our fixed-point operator $\mathbf{Ir}_i(\mathbf{A})$ faithfully captures the set in (5). The alternating reachability in the semantics implies that if $\mathbf{Ir}_i(\mathbf{A})$ holds at a world w and if $wR_i u$, then A_i and $\mathbf{Ir}_j(\mathbf{A})$ hold at world u , which corresponds to Lemma 2.2. Moreover, if $uR_i v$, then $\mathbf{Ir}_i(\mathbf{A})$ holds at world v , which corresponds to IRA_i . These reflect the self-referential structure shared by $\mathbf{Ir}_i(\mathbf{A})$ and $\mathbf{Ir}_j(\mathbf{A})$.

In addition, the proof of the above theorem gives the (strong) *finite model property* (cf., p.145, 339, Blackburn, *et al.* [1]). Thus, this logic is effectively decidable (called simply “decidable” in the logic literature), i.e., the set of provable formulae is recursive. In Section 6, we will discuss this problem relative to the game theoretic decidability/undecidability result for prediction/decision making.

The following lemma requires KD^2 to be the base logic for EIR^2 . It is proved by Theorem 2.1. If we add any of Axioms T, 4 or 5 to EIR^2 , the lemma does not hold. Counterexamples are given in Hu-Kaneko [7]. The failure of the following lemma under Axiom T is due to inseparability between player i 's mind and objective situation, which violates our basic approach to model player's subjective decision making in this paper.

Lemma 2.5. (Change of Scopes) (1): $\mathbf{B}_i(\Gamma_i^o) \vdash \mathbf{B}_i(A) \iff \Gamma_i^o \vdash A$;

(2): $\mathbf{B}_i(\Gamma_i^o) \vdash \neg \mathbf{B}_i(A) \iff \mathbf{B}_i(\Gamma_i^o) \vdash \mathbf{B}_i(\neg A)$.

In our applications, $\mathbf{B}_i(\Gamma_i^o)$ takes the form $\mathbf{Ir}_i(\mathbf{C})$, i.e., $\mathbf{Ir}_i(\mathbf{C}) \vdash \mathbf{B}_i(A)$ or $\mathbf{Ir}_i(\mathbf{C}) \vdash \neg \mathbf{B}_i(A)$. By Lemmas 2.2 and 2.5, this is equivalent to $\mathbf{Ir}_i^o(\mathbf{C}) \vdash A$ or $\mathbf{Ir}_i^o(\mathbf{C}) \vdash \neg A$. This is interpreted as meaning that $\mathbf{Ir}_i^o(\mathbf{C}) \vdash A$ or $\mathbf{Ir}_i^o(\mathbf{C}) \vdash \neg A$ is obtained in the mind of player i .

3 Game Theoretic Concepts

First, we give a few game theoretic concepts relevant for our discussions. Then, we formulate them in the language of EIR^2 . We also prepare some completeness results for game formulae, which are crucial to understand our game theoretic undecidability result.

3.1 Preliminary definitions

Let $G = (\{1, 2\}, \{S_1, S_2\}, \{h_1, h_2\})$ be a finite 2-person game, where $\{1, 2\}$ is the set of players, $S = S_1 \times S_2$ is the set of *strategy pairs*, and $h_i : S \rightarrow \mathbb{R}$ is the payoff function for player $i = 1, 2$. We write $(s_i; s_j)$ for $s = (s_1, s_2) \in S$. A strategy s_i for player i is a *best-response* against s_j iff $h_i(s_i; s_j) \geq h_i(t_i; s_j)$ for all $t_i \in S_i$. A strategy pair $s = (s_i; s_j)$ is a *Nash equilibrium* in G iff s_i is a best response against s_j for $i = 1, 2$. We denote the set of all Nash equilibria in G by $E(G)$. The set $E(G)$ may be empty, e.g., Table 1.3 has the empty $E(G)$. We say that s_i is a *Nash strategy* iff $(s_i; s_j)$ is a Nash equilibrium for some $s_j \in S_j$.

A subset E of S is *interchangeable* (Nash [16]) iff

$$\text{for all } s, s' \in E, (s_i; s'_j) \in E \text{ for } i = 1, 2. \quad (9)$$

This is equivalent to $E = E_1 \times E_2$, where $E_i = \{s_i \in S_i : (s_i; s_j) \in E \text{ for some } s_j\}$ for $i = 1, 2$. Let $\mathbf{E} = \{E : E \subseteq E(G) \text{ and } E \text{ satisfies (9)}\}$. The game G is *solvable* iff $E(G)$ satisfies (9), and we call $E(G)$ the Nash solution. Otherwise, it is *unsolvable*, and a nonempty set $F \subseteq S$ is a *subsolution* iff F is a maximal set in \mathbf{E} , i.e., there is no $E' \in \mathbf{E}$ such that $F \subsetneq E'$. Table 1.1 is solvable with the solution $\{(\mathbf{s}_{12}, \mathbf{s}_{21})\}$. Table 1.2 is unsolvable, and has two subsolutions: $\{(\mathbf{s}_{11}, \mathbf{s}_{21})\}$ and $\{(\mathbf{s}_{12}, \mathbf{s}_{22})\}$. Table 1.3 is solvable but has the empty $E(G)$ ⁸.

Hu-Kaneko [6] derived the Nash theory from the following decision criteria: Let E_i be a subset of S_i for $i = 1, 2$.

Na₁: for any $s_1 \in E_1$, s_1 is a best response against all $s_2 \in E_2$;

Na₂: for any $s_2 \in E_2$, s_2 is a best response against all $s_1 \in E_1$.

In Na_i , E_i describes the set of possible final decisions for player i , and E_j does i 's prediction about j 's possible final decisions. Here i 's prediction comes from his thinking about j 's criterion Na_j . When i makes his prediction based on Na_j , elements in E_j occur in the scope of j 's thinking, and this prediction occurs in the scope of i 's thinking. However, this argument is entirely interpretational. To make it explicit, we need the logic EIR².

The following proposition was proved in Hu-Kaneko [6].

Proposition 3.1. *Let $E(G) \neq \emptyset$, and E_i a nonempty subset of S_i for $i = 1, 2$.*

(1) *Suppose that G is solvable. Then $E = E_1 \times E_2$ is the Nash solution of G if and only if (E_1, E_2) is the greatest pair satisfying Na_1 - Na_2 .*⁹

(2) *Suppose that G is unsolvable. Then $E = E_1 \times E_2$ is a Nash subsolution if and only if (E_1, E_2) is a maximal pair satisfying Na_1 - Na_2 .*

These two cases correspond basically to the game theoretic decidability and undecidability results given in the subsequent sections. Here, we avoided unnecessary complication for the case of $E(G) = \emptyset$. In the subsequent sections, we treat this case, too.

⁸Nash [16] himself assumed the mixed strategies, and proved the existence of a Nash equilibrium. Here, we do not allow mixed strategies, and some games have no Nash equilibria.

⁹The "greatest" and "maximal" are relative to the componentwise set-inclusions.

3.2 Some completeness for game formulae

To express a game $G = (\{1, 2\}, \{S_1, S_2\}, \{h_1, h_2\})$ in EIR^2 , we formalize payoff functions h_1 and h_2 in terms of preference formulae (the players and strategies are already included in the language):

$$g_i = \wedge [\{\text{Pr}_i(s; t) : h_i(s) \geq h_i(t)\} \cup \{\neg\text{Pr}_i(s; t) : h_i(s) < h_i(t)\}]. \quad (10)$$

We call g_i the *formalized payoffs* associated with h_i for $i = 1, 2$. Here, $g = (g_1, g_2)$ is determined by G . Since (10) also contains negative preferences, for all $s, t \in S$, $g_i \vdash \text{Pr}_i(s; t)$ or $g_i \vdash \neg\text{Pr}_i(s; t)$, i.e., under g_i , completeness for all atomic preference formulae is obtained.

Consistency of $g_1 \wedge g_2$ can be shown by constructing a truth assignment. Consistency of the infinite regress $\mathbf{Ir}_i(g_1, g_2)$ in EIR^2 is also obtained by applying the belief eraser ε_0 : Suppose that $\mathbf{Ir}_i(g_1, g_2) \vdash \neg A \wedge A$ for some nonepistemic formula A . Applying ε_0 , we have $g_1 \wedge g_2 \vdash_0 \neg A \wedge A$ by (8), which is impossible because of consistency of $g_1 \wedge g_2$. In the same way, we have consistency of $\mathbf{Ir}_i^o(g_1, g_2)$ in EIR^2 . These are listed for reference:

$$\mathbf{Ir}_i(g_1, g_2) \text{ and } \mathbf{Ir}_i^o(g_1, g_2) \text{ are consistent in } \text{EIR}^2. \quad (11)$$

We formalize *best response* and *Nash equilibrium*: The statement “ $s_i \in S_i$ is a best response to $s_j \in S_j$ ” is expressed as $\text{bst}_i(s_i; s_j) := \wedge_{t_i \in S_i} \text{Pr}_i((s_i; s_j); (t_i; s_j))$. The statement “ $s = (s_1, s_2) \in S$ is a Nash equilibrium” is given as $\text{nash}(s) := \text{bst}_1(s_1; s_2) \wedge \text{bst}_2(s_2; s_1)$. The formulae defined above are game formulae.

Game theoretic undecidability could be an easy conclusion if a belief set for player i has a weak content. Thus, we assume that player i has enough beliefs, in order for our question to make sense. As far as game formulae are concerned, the infinite regress of the formalized payoffs $\mathbf{Ir}_i(g_1, g_2)$ contains sufficient information to prove or to disprove them.

Lemma 3.1. *Let A_i be a nonepistemic game formula for $i = 1, 2$. Let G be a game and $\mathbf{g} = (g_1, g_2)$ its formalized payoffs. Then,*

(1) $g_i \vdash A_i$ or $g_i \vdash \neg A_i$ for $i = 1, 2$;

(2) the following three are equivalent:

(a) $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{Ir}_i(\mathbf{A})$ for $i = 1, 2$; (b) $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_i^o(\mathbf{A})$ for $i = 1, 2$; (c) $g_i \vdash A_i$ for $i = 1, 2$.

Proof. (1) Let $\text{Pr}_i(s; t)$ be any atomic formula. Recall that $g_i \vdash \text{Pr}_i(s; t)$ or $g_i \vdash \neg\text{Pr}_i(s; t)$. We can extend this result to other nonepistemic game formulae for i by induction on their lengths.

(2) ((c) \implies (a) \implies (b)): Suppose that $g_i \vdash A_i$, i.e., $\vdash g_i \supset A_i$ for $i = 1, 2$. It follows from Lemma 2.3.(1) that $\vdash \mathbf{Ir}_i(g_1 \supset A_1, g_2 \supset A_2)$. By Lemma 2.3.(4), $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{Ir}_i(\mathbf{A})$ for $i = 1, 2$. Since $\vdash g_i \supset A_i$, we have $g_i \wedge \mathbf{Ir}_j(\mathbf{g}) \vdash A_i \wedge \mathbf{Ir}_j(\mathbf{A})$, i.e., $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_i^o(\mathbf{A})$.

((b) \implies (c)): Suppose that $g_1 \not\vdash A_1$ or $g_2 \not\vdash A_2$. By (1), $g_i \vdash \neg A_i$ or $g_j \vdash \neg A_j$ or both. We only consider the case where $g_i \vdash A_i$ and $g_j \vdash \neg A_j$. Using the same arguments as above, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_i^o(A_i; \neg A_j)$. By Lemma 2.4.(1), $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{B}_j(\neg A_j)$ and hence, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{B}_j(A_j)$. But by Lemma 2.4.(1), $\vdash \mathbf{Ir}_i^o(\mathbf{A}) \supset \mathbf{B}_j(A_j)$, equivalently, $\vdash \neg \mathbf{B}_j(A_j) \supset \neg \mathbf{Ir}_i^o(\mathbf{A})$. Thus, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{Ir}_i^o(A_i; A_j)$. By (11), we have $\mathbf{Ir}_i^o(\mathbf{g}) \not\vdash \mathbf{Ir}_i^o(A_i; A_j)$. The other cases are similar. ■

The next theorem shows that $\mathbf{Ir}_i(\mathbf{g})$ is complete relative to infinite regresses of nonepistemic game formulae. It states this in terms of the epistemic content $\mathbf{Ir}_i^o(\cdot; \cdot)$ for coherency of the later purpose.

Theorem 3.1. (Completeness for infinite regresses of game formulae) Let G be a game and $\mathbf{g} = (g_1, g_2)$ its formalized payoffs. Let A_i be a nonepistemic game formula for $i = 1, 2$. Then, either $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_i^o(\mathbf{A})$ or $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{Ir}_i^o(\mathbf{A})$, which implies either $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{Ir}_i(\mathbf{A})$ or $\mathbf{Ir}_i(\mathbf{g}) \vdash \neg \mathbf{Ir}_i(\mathbf{A})$.

Proof. Since $g_i \vdash A_i$ or $g_i \vdash \neg A_i$ for $i = 1, 2$, we should consider the four cases. Here, we consider only the case where $g_i \vdash \neg A_i$ for $i = 1, 2$. By (6), $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg A_i$. Using the contrapositive of Lemma 2.4.(1), we have $\vdash \neg A_i \supset \neg \mathbf{Ir}_i^o(A_i; A_i)$. Thus, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{Ir}_i^o(A_i; A_i)$. ■

The above theorem, together with the next one, will be used for our game theoretic decidability result. In fact, the result gets sharper with Axiom T, i.e., $\text{EIR}^2(\text{T})$. In particular, the next theorem will be used for the full completeness (Theorem 4.4) for solvable games and the no-formula theorem (Theorem 5.4) for unsolvable games.

Theorem 3.2. (Completeness for game formulae under Axiom T) Let G be a game and $\mathbf{g} = (g_1, g_2)$ its formalized payoffs. For any game formula A , either $\mathbf{Ir}_i(\mathbf{g}) \vdash A$ or $\mathbf{Ir}_i(\mathbf{g}) \vdash \neg A$ in $\text{EIR}^2(\text{T})$.

Proof. We prove the claim $\mathbf{Ir}_i^o(\mathbf{g}) \vdash A$ or $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg A$ by induction on the length of A . This implies $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{B}_i(A)$ or $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg A)$; then we have the assertion by Axiom T. Let A be an atomic formula. Then, $g_1 \wedge g_2 \vdash A$ or $g_1 \wedge g_2 \vdash \neg A$. Then, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash g_1 \wedge g_2$ by (6) and Axiom T. Thus, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash A$ or $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg A$.

Let A be nonatomic, and suppose the inductive hypothesis that decidability holds for the immediate subformulae of A . Let $A = C \supset D$. By the inductive hypothesis, decidability holds for C and D . Using this, we have $\mathbf{Ir}_i^o(\mathbf{g}) \vdash A$ or $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg A$. Similar arguments apply to connectives \wedge, \vee and \neg .

Let $A = \mathbf{B}_k(C)$. The hypothesis is: $\mathbf{Ir}_i^o(\mathbf{g}) \vdash C$ or $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg C$. Let $\mathbf{Ir}_i^o(\mathbf{g}) \vdash C$. Then, $\mathbf{B}_k(\mathbf{Ir}_i^o(\mathbf{g})) \vdash \mathbf{B}_k(C)$. By IRA_i^o and Axiom T, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{B}_j(\mathbf{Ir}_i^o(\mathbf{g}))$ and $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{Ir}_i^o(\mathbf{g}))$. Thus, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{B}_k(C)$. Now, let $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg C$. By the same arguments, we have $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{B}_k(\neg C)$, and, by Axiom D, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{B}_k(C)$.

Let $A = \mathbf{Ir}_k(C_1, C_2)$. The induction hypothesis is that decidability holds for C_1 and C_2 . Now, suppose $\mathbf{Ir}_i^o(\mathbf{g}) \vdash C_1 \wedge C_2$. As remarked for $\text{EIR}^2(\text{T})$ in the end of Section 2.2, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_j^o(\mathbf{g})$ and $\mathbf{Ir}_j^o(\mathbf{g}) \vdash \mathbf{Ir}_i^o(\mathbf{g})$. Hence, $\mathbf{Ir}_k^o(\mathbf{g}) \vdash C_k$ for $k = 1, 2$. Thus, $\mathbf{Ir}_k(\mathbf{g}) \vdash \mathbf{B}_k(C_k)$ for $k = 1, 2$. By Lemma 2.3 (1), $\mathbf{Ir}_k(\mathbf{g}) \vdash \mathbf{Ir}_k(C_1, C_2)$ for $k = 1, 2$. Since $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_k(\mathbf{g})$ for $k = 1, 2$ by (6) and Axiom T, we have $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_k(C_1, C_2)$.

Let $\mathbf{Ir}_i^o(\mathbf{g}) \vdash (\neg C_i) \wedge C_j$. By the same argument, we have $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \mathbf{Ir}_i(\neg C_i; C_j)$. By Lemma 2.3.(5), $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg \mathbf{Ir}_i(C_i; C_j)$. The same argument can be applied to the case of $\mathbf{Ir}_i^o(\mathbf{g}) \vdash C_i \wedge (\neg C_j)$ and $\mathbf{Ir}_i^o(\mathbf{g}) \vdash (\neg C_i) \wedge (\neg C_j)$. ■

4 Formalized Nash Theory

We give three axioms for player i 's prediction/decision making, and assume the symmetric axioms for player i 's prediction about player j 's prediction/decision making. These lead to an infinite regress of those axioms. In this section, we show, for a solvable game, that the infinite regress of those axioms can be fully explicated, and obtain the decidability result.

4.1 Axioms for Prediction/Decision Making

We start with the following three axioms. These are described in the mind of player i , i.e., in the scope of $\mathbf{B}_i(\cdot)$:

N0_{*i*} (Optimization against all predictions): $\bigwedge_{s \in S} [I_i(s_i) \wedge \mathbf{B}_j(I_j(s_j)) \supset \text{bst}_i(s_i; s_j)]$.

N1_{*i*} (Necessity of predictions): $\bigwedge_{s_i \in S_i} \langle I_i(s_i) \supset \bigvee_{s_j \in S_j} \mathbf{B}_j(I_j(s_j)) \rangle$.

N2_{*i*} (Predictability): $\bigwedge_{s_i \in S_i} \langle I_i(s_i) \supset \mathbf{B}_j \mathbf{B}_i(I_i(s_i)) \rangle$.

For each $i = 1, 2$, let $\mathbf{N}_i = \mathbf{N0}_i \wedge \mathbf{N1}_i \wedge \mathbf{N2}_i$, and let $\mathbf{N} = (\mathbf{N}_1, \mathbf{N}_2)$.

The first axiom directly corresponds to Na_i . The second requires player i to have a prediction for his decision. It corresponds to the nonemptiness of E_1 and E_2 in Proposition 3.1, while $\mathbf{N1}_i$ allows both to be empty. The third states that in the mind of player i , his decision is correctly predicted by player j . We find a similar structure in Axiom IRA_i , but note that $\mathbf{N2}_i$ and IRA_i have different orders of applications of \mathbf{B}_i and \mathbf{B}_j . Indeed, $I_i(s_i)$ is rather naked without having the intended scope of $\mathbf{B}_i(\cdot)$, while $\mathbf{I}_i(\cdot, \cdot)$ includes the outer $\mathbf{B}_i(\cdot)$, shown as in Lemma 2.2.

Axioms \mathbf{N}_i and \mathbf{N}_j are interdependent: \mathbf{N}_i is assumed in the mind of player i , i.e., $\mathbf{B}_i(\mathbf{N}_i)$. Since \mathbf{N}_i includes $\mathbf{B}_j(I_j(s_j))$, player i needs to predict what j would choose. This prediction is made by the criterion $\mathbf{B}_i \mathbf{B}_j(\mathbf{N}_j)$. Then, $\mathbf{B}_i(I_i(s_i))$ requires $\mathbf{B}_i \mathbf{B}_j \mathbf{B}_i(\mathbf{N}_i)$, and so on. These are captured by the infinite regress formula $\mathbf{I}_i(\mathbf{N}) = \mathbf{I}_i(\mathbf{N}_i; \mathbf{N}_j)$. The infinite regress $\mathbf{I}_i(\mathbf{N})$ within the logic EIR^2 may be compared with Johansen's [9] interpretation of Nash theory. This will be discussed in Section 6.

We take the infinite regress $\mathbf{I}_i(\mathbf{N}_i; \mathbf{N}_j)$ as basic beliefs for player i 's prediction/decision making; $I_i(s_i)$ and $\mathbf{B}_j(I_j(s_j))$ in $\mathbf{I}_i(\mathbf{N}_i; \mathbf{N}_j)$ are treated as "unknowns" to be found by player i with logical analysis. From $\mathbf{I}_i(\mathbf{N}_i; \mathbf{N}_j)$, necessary conditions for $I_i(s_i)$ and $I_j(s_j)$ are derived as the following game formulae: for each $i = 1, 2$ and $s_i \in S_i$,

$$A_i^*(s_i) := \bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)]. \quad (12)$$

These candidate formulae play a crucial role in our subsequent analysis.

The nonepistemic content of $A_i^*(s_i)$ is given as $\varepsilon_0(A_i^*(s_i)) = \bigvee_{t_j \in S_j} \langle \text{bst}_i(s_i; t_j) \wedge \text{bst}_j(t_j; s_i) \rangle = \bigvee_{t_j \in S_j} \text{nash}(s_i; t_j)$. That is, $\varepsilon_0(A_i^*(s_i))$ means " s_i is a Nash strategy". In the logic $\text{EIR}^2(\mathbf{T})$, we may interpret $\mathbf{I}_i(\cdot, \cdot)$ as the common knowledge operator (recall cka and cki in Section 2.2), and hence $A_i^*(s_i)$ means " s_i is a common knowledge Nash strategy". We emphasize this interpretation with Axiom \mathbf{T} by writing $\mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)]$ as $\mathbf{C}^*(\text{Nash}(s_i; t_j))$ in $\text{EIR}^2(\mathbf{T})$, and $A_i^*(s_i)$ becomes $\bigvee_{t_j \in S_j} \mathbf{C}^*(\text{Nash}(s_i; t_j))$. This formula was discussed in Kaneko-Nagashima [10] and Kaneko [12]. While without Axiom \mathbf{T} , the formula $A_i^*(s_i)$ occurs in the mind of player i , independent of reality as well as player j , with Axiom \mathbf{T} , $\bigvee_{t_j \in S_j} \mathbf{C}^*(\text{Nash}(s_i; t_j))$ ($\equiv A_i^*(s_i)$) describes reality as well as both players' thinking.

We have the following result, which will be proved in the end of this subsection.

Theorem 4.1. (Necessity) For $i = 1, 2$,

$$\mathbf{I}_i(\mathbf{N}) \vdash \mathbf{B}_i(I_i(s_i) \supset A_i^*(s_i)) \text{ for all } s_i \in S_i. \quad (13)$$

That is, player i infers $A_i^*(s_i)$ as a necessary condition for his decision. By this and Lemma 2.2, we have also $\mathbf{I}_i(\mathbf{N}) \vdash \mathbf{B}_i[\mathbf{B}_j(I_j(s_j)) \supset \mathbf{B}_j(A_j^*(s_j))]$ for all $s_j \in S_j$; player i infers $\mathbf{B}_j(A_j^*(s_j))$ as a necessary conditions for his prediction. By Lemma 2.3.(1), we have, also,

$\mathbf{Ir}_i(\mathbf{N}) \vdash \mathbf{Ir}_i[\mathbf{I}_i(s_i) \supset A_i^*(s_i); \mathbf{I}_j(s_j) \supset A_j^*(s_j)]$ for all $s \in S$. That is, those necessary conditions form an infinite regress, too. From now on, we focus on statements of the form of (13).

Recalling $\varepsilon_0(A_i^*(s_i))$, (13) may be interpreted as meaning that a Nash strategy is derived. However, our target is prediction/decision making by a player. A possible decision resulting from this process is expressed by $\mathbf{I}_i(s_i)$, and $A_i^*(s_i)$ is only a necessary condition for it as a purely solution-theoretic statement without specifying payoffs. Also, even if payoffs, e.g., $\mathbf{Ir}_i(g_1, g_2)$, are specified, (13) does not give a positive answer to $\mathbf{I}_i(s_i)$; that is, the contrapositive of (13) may give only a negative decision $\neg \mathbf{I}_i(s_i)$ from $\neg A_i^*(s_i)$. We discuss the converse of (13) under the assumption of $\mathbf{Ir}_i(g_1, g_2)$ in later sections.

Here, we prove Theorem 4.1. It follows from (2) of the next lemma. (1) does not need $\mathbf{N}1_i$. We write $\mathbf{N}0_i \wedge \mathbf{N}2_i$ as $\mathbf{N}02_i$ for $i = 1, 2$.

Lemma 4.1. *For $i = 1, 2$, and $s = (s_i; s_j) \in S$,*

(1): $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \vdash \mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset \mathbf{Ir}_i^\circ[\mathbf{bst}_i(s_i; s_j); \mathbf{bst}_j(s_i; s_j)]$;

(2): $\mathbf{Ir}_i^\circ[\mathbf{N}_i; \mathbf{N}_j] \vdash \mathbf{I}_i(s_i) \supset A_i^*(s_i)$.

Proof. (1): Let $\theta_i(s_i; s_j) := \mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \wedge \mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j))$. Here, we show, for $i = 1, 2$,

$$\vdash \theta_i(s_i; s_j) \supset \mathbf{bst}_i(s_i; s_j) \wedge \mathbf{B}_j(\mathbf{bst}_j(s_j; s_i)) \wedge \mathbf{B}_j \mathbf{B}_i(\theta_i(s_i, s_j)). \quad (14)$$

By this and Lemma 2.4.(2), we have $\vdash \theta_i(s_i; s_j) \supset \mathbf{Ir}_i^\circ[\mathbf{bst}_i(s_i; s_j); \mathbf{bst}_j(s_i; s_j)]$, which implies the assertion.

The first part, $\vdash \theta_i(s_i; s_j) \supset \mathbf{bst}_i(s_i; s_j)$, of (14) comes from $\mathbf{N}0_i$ and $\mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j))$. Consider the second part. Since $\vdash \theta_i(s_i; s_j) \supset \mathbf{B}_j(\mathbf{N}02_j)$ and $\vdash \mathbf{B}_j(\mathbf{N}02_j) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \wedge \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i(s_i)) \supset \mathbf{B}_j(\mathbf{bst}_j(s_j; s_i))$, we have $\vdash \theta_i(s_i; s_j) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \wedge \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i(s_i)) \supset \mathbf{B}_j(\mathbf{bst}_j(s_j; s_i))$. Observe that $\mathbf{B}_j(\mathbf{I}_j(s_j))$ is included in $\theta_i(s_i, s_j)$ and $\mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i(s_i))$ is derived from $\mathbf{I}_i(s_i)$ in $\theta_i(s_i; s_j)$ by $\mathbf{N}2_i$. Hence, $\vdash \theta_i(s_i; s_j) \supset \mathbf{B}_j(\mathbf{bst}_j(s_j; s_i))$. Now, consider the third part of (14). By Lemma 2.4.(1), $\vdash \mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \supset \mathbf{B}_j \mathbf{B}_i(\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j])$. Using $\mathbf{N}2_i$, we have $\vdash \mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \wedge \mathbf{I}_i(s_i) \supset \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i(s_i))$, and, using $\mathbf{B}_j(\mathbf{N}2_j)$ in $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j]$, we have $\vdash \mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset \mathbf{B}_j \mathbf{B}_i \mathbf{B}_j(\mathbf{I}_j(s_j))$. Summing those three up, we obtain $\vdash \theta_i(s_i; s_j) \supset \mathbf{B}_j \mathbf{B}_i(\theta_i(s_i; s_j))$.

(2): It follows from (1) that $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \vdash \mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset \bigvee_{t_j \in S_j} \mathbf{Ir}_i^\circ[\mathbf{bst}_i(s_i; t_j); \mathbf{bst}_j(t_j; s_i)]$. This is equivalent to $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \vdash \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset (\mathbf{I}_i(s_i) \supset A_i^*(s_i))$. Hence $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j] \vdash \bigvee_{t_j \in S_j} \mathbf{B}_j(\mathbf{I}_j(t_j)) \supset (\mathbf{I}_i(s_i) \supset A_i^*(s_i))$. Adding $\mathbf{N}1_i$ to $\mathbf{Ir}_i^\circ[\mathbf{N}02_i; \mathbf{N}02_j]$, we delete the first disjunctive formula, i.e., $\mathbf{Ir}_i^\circ[\mathbf{N}_i; \mathbf{N}_j] \vdash \mathbf{I}_i(s_i) \supset A_i^*(s_i)$. ■

4.2 Choice of the deductive weakest formulae for \mathbf{N}_i and \mathbf{N}_j

The basic belief $\mathbf{Ir}_i[\mathbf{N}_i; \mathbf{N}_j]$ only gives necessary conditions for $\mathbf{I}_i(s_i)$ and $\mathbf{B}_j(\mathbf{I}_j(s_j))$, but not sufficient conditions. In fact, there are other formulae than $A_i^*(s_i)$ and $A_j^*(s_j)$ that enjoy the properties described by \mathbf{N}_i and \mathbf{N}_j . For example, the families of formulae, $\{\perp(s_i)\}_{s_i \in S_i}$, $i = 1, 2$, where $\perp(s_i) := \neg(p \supset p)$, $s_i \in S_i$ and p is a atomic preference formula, makes $\mathbf{N}_i = \mathbf{N}0_i \wedge \mathbf{N}1_i \wedge \mathbf{N}2_i$ trivially hold with the substitution of $\perp(s_i)$ for each $\mathbf{I}_i(s_i)$ in \mathbf{N}_i . To avoid such unintended candidates and to analyze the exact logical contents of $\mathbf{Ir}_i[\mathbf{N}_i; \mathbf{N}_j]$, we choose families of formulae $\{A_i(s_i)\}_{s_i \in S_i}$ and $\{A_j(s_j)\}_{s_j \in S_j}$ having *only* the properties \mathbf{N}_i and \mathbf{N}_j .

We formalize this choice by an axiom scheme. Let $\mathcal{A} = (\mathcal{A}_i; \mathcal{A}_j)$ be a pair of *candidate families* $\mathcal{A}_i = \{A_i(s_i)\}_{s_i \in S_i}$ and $\mathcal{A}_j = \{A_j(s_j)\}_{s_j \in S_j}$. Let $\mathbf{N}_i(\mathcal{A})$ be the formula obtained from

N_i by substituting $(A_1(s_1), A_2(s_2))$ for $(I_1(s_1), I_2(s_2))$ for each $s = (s_1, s_2) \in S$. We denote the following formula by $WF_i(\mathcal{A})$:

$$\begin{aligned} N_i(\mathcal{A}) \wedge \mathbf{B}_j(N_j(\mathcal{A})) \wedge [\wedge_{s \in S} \langle I_i(s_i) \wedge \mathbf{B}_j(I_j(s_j)) \supset A_i(s_i) \wedge \mathbf{B}_j(A_j(s_j)) \rangle] & \quad (15) \\ & \supset \wedge_{s_i \in S_i} \langle A_i(s_i) \supset I_i(s_i) \rangle. \end{aligned}$$

Let $\mathbf{WF}(\mathcal{A}) = (WF_1(\mathcal{A}), WF_2(\mathcal{A}))$. The axiom scheme for the choice of the weakest candidate formulae is denoted by $\mathbf{Ir}_i(\mathbf{WF})$, i.e., it is the set $\{\mathbf{Ir}_i(\mathbf{WF}(\mathcal{A})) : \mathcal{A} \text{ is a pair of candidate families}\}$.

The formula $WF_i(\mathcal{A})$ in (15) contains the additional premise $\wedge_{s \in S} \langle I_i(s_i) \wedge \mathbf{B}_j(I_j(s_j)) \supset A_i(s_i) \wedge \mathbf{B}_j(A_j(s_j)) \rangle$. A sole use of $WF_i(\mathcal{A})$ is not meaningful since $I_i(s_i) \wedge \mathbf{B}_j(I_j(s_j))$ have no properties, yet. It is used together with $\mathbf{Ir}_i(N_i; N_j)$. This premise corresponds to the maximality requirement in the definition of a subsolution in Section 3. If we drop the additional premise, (15) becomes

$$WF_i^+(\mathcal{A}) = N_i(\mathcal{A}) \wedge \mathbf{B}_j(N_j(\mathcal{A})) \supset \wedge_{s_i \in S_i} \{A_i(s_i) \supset I_i(s_i)\}. \quad (16)$$

This is stronger than $WF_i(\mathcal{A})$. This $WF_i^+(\mathcal{A})$ works only for a solvable game, but not for an unsolvable game, while $WF_i(\mathcal{A})$ in (15) works for any game.

We study implications from $\{\mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF})$ under the infinite regress of formalized payoffs $\mathbf{Ir}_i(\mathbf{g}) = \mathbf{Ir}_i(g_i; g_j)$. We postulate the entire set of axioms, denoted by $\Delta_i(\mathbf{g}) := \{\mathbf{Ir}_i(\mathbf{g}), \mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF})$, as the basic beliefs for player i 's prediction/decision making.

We first state the consistency of the basic beliefs $\Delta_i(\mathbf{g})$. The following lemma will be proved in the proof of Lemma 5.1.

Lemma 4.2. (*Consistency of the belief set*) $\Delta_i(\mathbf{g})$ is consistent for any game G .

In fact, $\Delta_i^+(\mathbf{g}) = \{\mathbf{Ir}_i(\mathbf{g}), \mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF}^+)$ is consistent if and only if G is a solvable game, and Δ_i^+ is equivalent to Δ_i for any solvable G .

The formalized Nash theory is expressed as $(\text{EIR}^2; \Delta_i(\mathbf{g}))$. That is, we fix the logical system EIR^2 , and within it, we have the set of nonlogical axioms $\Delta_i(\mathbf{g})$, which depends upon a game G . We are interested in the logical implications related to prediction/decision making derived from $\Delta_i(\mathbf{g})$ in EIR^2 .

4.3 Game theoretic decidability for solvable games

Here, we show that the basic beliefs $\Delta_i(\mathbf{g})$ determine the possible final decisions for a solvable game. The proof of this theorem is given in the end of this subsection.

Theorem 4.2. (*Determination I*) Let G be a solvable game and \mathbf{g} its formalized payoffs. Then, for $i = 1, 2, W$

$$\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(I_i(s_i) \equiv A_i^*(s_i)) \text{ for all } s_i \in S_i. \quad (17)$$

Proof. We prove the following claims.

Claim 1: Let G be solvable. Then, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash A_i^*(s_i) \wedge \mathbf{B}_j(A_j^*(s_j)) \supset \text{bst}_i(s_i; s_j)$.

Claim 2: $\vdash A_i^*(s_i) \supset \vee_{t_j \in S_j} \mathbf{B}_j(A_j^*(t_j))$.

Claim 3: $\vdash A_i^*(s_i) \supset \mathbf{B}_j \mathbf{B}_i(A_i^*(s_i))$.

Proof of Claim 1: Since $\text{bst}_i(s_i; s_j)$ is a game formula for $i = 1, 2$, we have, for each $s \in S$, $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{I}_i^o(\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i))$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg \mathbf{I}_i^o(\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i))$ by Theorem 3.1. Hence, for each $s_i \in S_i$, $\mathbf{I}_i^o(\mathbf{g}) \vdash A_i^*(s_i)$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg A_i^*(s_i)$. Using Lemma 2.2, we have, for each $s_j \in S_j$, $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{B}_j(A_j^*(s_j))$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg \mathbf{B}_j(A_j^*(s_j))$. Also, for each $s \in S$, $\mathbf{I}_i^o(\mathbf{g}) \vdash \text{bst}_i(s_i; s_j)$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg \text{bst}_i(s_i; s_j)$. Thus, $\mathbf{I}_i^o(\mathbf{g}) \vdash A_i^*(s_i) \wedge \mathbf{B}_j(A_j^*(s_j)) \supset \text{bst}_i(s_i; s_j)$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg[A_i^*(s_i) \wedge \mathbf{B}_j(A_j^*(s_j)) \supset \text{bst}_i(s_i; s_j)]$. If the latter held, then, applying the epistemic eraser ε_0 to this, we would have $g_i \wedge g_j \vdash \neg[(\bigvee_{t_j \in S_j} \text{nash}(s_i, t_j)) \wedge (\bigvee_{t_i \in S_i} \text{nash}(s_j, t_i)) \supset \text{bst}_i(s_i; s_j)]$, which is impossible since G is a solvable game. Hence, we have the assertion.

Proof of Claim 2: By Lemma 2.2, we have $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \mathbf{B}_j(\mathbf{I}_j^o[\text{bst}_j(s_j; s_i); \text{bst}_i(s_i; s_j)])$. Hence, $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \mathbf{B}_j(\bigvee_{t_i \in S_i} \mathbf{I}_j^o[\text{bst}_j(s_j; t_i); \text{bst}_i(s_i; t_j)])$, i.e., $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \mathbf{B}_j(A_j^*(s_j))$. Hence, $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(A_j^*(t_j))$. Then, $\vdash \bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)] \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(A_j^*(t_j))$, i.e., $\vdash A_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(A_j^*(t_j))$.

Proof of Claim 3: Since $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \mathbf{B}_j(\mathbf{I}_j^o[\text{bst}_j(s_j; s_i); \text{bst}_i(s_i; s_j)])$ and $\vdash \mathbf{B}_j(\mathbf{I}_j^o[\text{bst}_j(s_j; s_i); \text{bst}_i(s_i; s_j)]) \supset \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)])$, we have $\vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)] \supset \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)])$. We take disjunctions from the latter to the former with respect to s_j , and have $\vdash \bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)] \supset \bigvee_{t_j \in S_j} \mathbf{B}_j \mathbf{B}_i(\mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)])$. Then, the former is $A_i^*(s_i)$, and the latter implies $\mathbf{B}_j \mathbf{B}_i(\bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)])$, i.e., $\mathbf{B}_j \mathbf{B}_i(A_i^*(s_i))$.

Here, we prove the theorem. It follows from the above claims that $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{N}_i(\mathcal{A}^*)$ for $i = 1, 2$. Hence, $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{N}_i(\mathcal{A}^*) \wedge \mathbf{B}_j(\mathbf{N}_j(\mathcal{A}^*))$. It follows from Theorem 4.1 that $\mathbf{I}_i^o(\mathbf{N}_i; \mathbf{N}_j) \vdash \bigwedge_{s \in S} [\mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset A_i^*(s_i) \wedge \mathbf{B}_j(A_j^*(s_j))]$. We have $\mathbf{I}_i^o(\mathbf{g}), \mathbf{I}_i^o(\mathbf{N}), \mathbf{I}_i^o(\mathbf{WF}) \vdash [A_i^*(s_i) \supset \mathbf{I}_i(s_i)] \wedge [\mathbf{B}_j(A_j^*(s_j)) \supset \mathbf{B}_j(\mathbf{I}_j(s_j))]$. Hence, $\mathbf{I}_i^o(\mathbf{g}), \mathbf{I}_i^o(\mathbf{N}), \mathbf{I}_i^o(\mathbf{WF}) \vdash \mathbf{I}_i^o[A_i^*(s_i) \supset \mathbf{I}_i(s_i); A_i^*(s_j) \supset \mathbf{I}_i(s_j)]$. Using Theorem 4.1 and the Claim (3), we have $\mathbf{I}_i(\mathbf{g}), \mathbf{I}_i(\mathbf{N}), \mathbf{I}_i(\mathbf{WF}) \vdash \mathbf{I}_i[A_i^*(s_i) \equiv \mathbf{I}_i(s_i); A_i^*(s_j) \equiv \mathbf{I}_i(s_j)]$. ■

That is, player i infers from his beliefs $\Delta_i(\mathbf{g})$ that his possible decision and prediction are fully expressed by $A_i^*(s_i)$ and $\mathbf{B}_j(A_j^*(s_j))$ for a solvable game G . As remarked above, in the logic $\text{EIR}^2(\mathbf{T})$, $A_i^*(s_i)$ can be written as $\bigvee_{t_j \in S_j} \mathbf{C}^*(\text{Nash}(s_i; t_j))$, and Theorem 4.2 becomes $\Delta_i(\mathbf{g}) \vdash \mathbf{I}_i(s_i) \equiv \bigvee_{t_j \in S_j} \mathbf{C}^*(\text{Nash}(s_i; t_j))$. That is, a possible decision s_i is the Nash strategy with common knowledge. This corresponds to the result given in Kaneko [11].

Then, because of the above theorem and Theorem 3.1, player i can decide whether a given strategy s_i is a final decision for him or not, which is stated by the following theorem.

Theorem 4.3. (Game theoretic decidability) *Let G be a solvable game and $\mathbf{g} = (g_1, g_2)$ its formalized payoffs. Then, for $i = 1, 2$ and each $s_i \in S_i$,*

$$\text{either } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i)) \text{ or } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg \mathbf{I}_i(s_i)). \quad (18)$$

Proof. We show (18). Since $\text{bst}_i(s_i; s_j)$ is a nonepistemic game formula for i , it follows from Theorem 3.1 that $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)]$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)]$. If s_i is a Nash strategy for G , then $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{I}_i^o[\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)]$ for some $s_j \in S_j$; so, $\mathbf{I}_i^o(\mathbf{g}) \vdash \bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)]$, i.e., $\mathbf{I}_i^o(\mathbf{g}) \vdash A_i^*(s_i)$. If not, we have $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg \bigvee_{t_j \in S_j} \mathbf{I}_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)]$, i.e., $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg A_i^*(s_i)$. Thus, we have $\mathbf{I}_i(\mathbf{g}) \vdash \mathbf{B}_i(A_i^*(s_i))$ or $\mathbf{I}_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg A_i^*(s_i))$. By (17), we have $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$ or $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg \mathbf{I}_i(s_i))$. ■

By Theorem 4.3 and Lemma 2.5, we also have, for each strategy $s_j \in S_j$,

$$\text{either } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i\mathbf{B}_j(I_j(s_j)) \text{ or } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i\mathbf{B}_j(\neg I_j(s_j)). \quad (19)$$

Thus, player i can predict whether a given strategy s_j for j is a possible decision for him for not. From now on, we concentrate on decidability or undecidability for player i .

Since $\varepsilon_0 A_i^*(s_i) = \bigvee_{t_j \in S_j} \text{nash}(s_i; t_j)$, the positive or negative decision in (18) corresponds to whether s_i is a Nash strategy or not. For the negative case, we need to add only $\mathbf{I}_i(\mathbf{g})$ to $\mathbf{I}_i(\mathbf{N})$ in Theorem 4.1, that is, if s_i is not a Nash strategy, then

$$\mathbf{I}_i(\mathbf{g}), \mathbf{I}_i(\mathbf{N}) \vdash \mathbf{B}_i(\neg I_i(s_i)). \quad (20)$$

This result is independent of the solvability of the game G . For the positive case, we need the full set $\Delta_i(\mathbf{g}) = \{\mathbf{I}_i(\mathbf{g}), \mathbf{I}_i(\mathbf{N})\} \cup \mathbf{I}_i(\mathbf{WF})$ and the solvability of G .

Since Table 1.1 is a solvable game, Theorem 4.3 is applicable, and the belief set $\Delta_1(\mathbf{g})$ recommends strategy \mathbf{s}_{12} as a positive decision to player 1, but $\mathbf{s}_{11}, \mathbf{s}_{13}$ as negative decisions. By (19), player 2 would choose \mathbf{s}_{21} , and would deny the others. Table 1.2 is an unsolvable game; Theorem 4.2 is not applicable. In Table 1.3, (20) recommends all strategies as negative decisions.

In the logic $\text{EIR}^2(\mathbf{T})$, Theorem 4.3 becomes the full completeness theorem: the following theorem states that the theory $(\text{EIR}^2(\mathbf{T}); \Delta_i(g))$ is complete. From the game theoretic perspective, Theorem 4.3 is sufficient for our purpose. However, at expense of the subjective nature for decision/prediction making, we obtain full completeness with Axiom T, which gives a full characterization of logical contents of $\Delta_i(g)$. Moreover, as a corollary, $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$ is effectively decidable.

Theorem 4.4. (Full Completeness with Axiom T) *Let G be a solvable game. Then, the theory $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$ is complete, i.e., for any $A \in \mathcal{P}$, $\Delta_i(\mathbf{g}) \vdash A$ or $\Delta_i(\mathbf{g}) \vdash \neg A$.*

Proof. It holds that $\Delta_i(\mathbf{g}) \vdash I_i(s_i) \equiv A_i^*(s_i)$ for any $s_i \in S_i$ and $i = 1, 2$ in $\text{EIR}^2(\mathbf{T})$. Let C be any formula, and $C^\#$ the formula obtained by replacing each occurrence of $I_i(s_i)$ in C by $A_i^*(s_i)$ ($s_i \in S_i, i = 1, 2$). We can show by induction of the length of a formula that $\Delta_i(\mathbf{g}) \vdash C^\# \equiv C$. We consider only the case of $C = \mathbf{I}_i(C_1, C_2)$. The induction hypothesis is that $\Delta_i(\mathbf{g}) \vdash C_k^\# \equiv C_k$ for $k = 1, 2$. Recall $\Delta_i(\mathbf{g}) \vdash A$ implies $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_k(A)$ for $k = 1, 2$ in $\text{EIR}^2(\mathbf{T})$. It follows from IRA_i that $\Delta_i(\mathbf{g}) \vdash \mathbf{I}_i(C_1, C_2) \supset \mathbf{B}_i(C_1^\#) \wedge \mathbf{B}_i\mathbf{B}_j(C_2^\#) \wedge \mathbf{B}_i\mathbf{B}_j(\mathbf{I}_i(C_1, C_2))$. By IRI_i , we have $\Delta_i(\mathbf{g}) \vdash \mathbf{I}_i(C_1, C_2) \supset \mathbf{I}_i(C_1^\#, C_2^\#)$. The converse is parallel. Then, it follows from Theorem 3.2 that $\Delta_i(\mathbf{g}) \vdash C$ or $\Delta_i(\mathbf{g}) \vdash \neg C$. ■

5 Game Theoretic Undecidability for Unsolvable Games

The situation for an unsolvable game differs entirely from that for a solvable game. When G is unsolvable, we have the undecidability result that for some strategy s_i for player i , he cannot infer from his belief set $\Delta_i(\mathbf{g}) = \{\mathbf{I}_i(\mathbf{g}), \mathbf{I}_i(\mathbf{N})\} \cup \mathbf{I}_i(\mathbf{WF})$ whether s_i is a final decision or not. We give three other results related to this theorem.

5.1 Game theoretic undecidability

Here is the main result of the paper. We place all the proofs of the results in this section in Section 5.2.

Theorem 5.1. (Game theoretic undecidability) Let G be an unsolvable game, $\mathbf{g} = (g_1, g_2)$ its formalized payoffs, and $i = 1, 2$. Then, there is an $s_i \in S_i$ such that

$$\text{neither } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(I_i(s_i)) \text{ nor } \Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg I_i(s_i)). \quad (21)$$

This result also holds in the logic $EIR^2(T)$.

This result differs from the negative result for a game with no Nash equilibria: For such a game, Theorem 4.3 states that player i can deny any strategy for his decision. In this case, he may think about some other criterion. In contrast, undecidability means that he can not reach such a conclusion.

However, the negative decision given in (20) holds for a non-Nash strategy s_i for any game G . Hence, s_i for (21) has to be a Nash strategy. Later, we show that a necessary and sufficient condition for (21) is that

$$s_i \text{ is a Nash strategy but } s_i \notin F_i \text{ for some subsolution } F_1 \times F_2. \quad (22)$$

Here, we give two examples. The battle of the sexes (Table 1.2) has two subsolutions $\{(\mathbf{s}_{11}, \mathbf{s}_{21})\}$, $\{(\mathbf{s}_{12}, \mathbf{s}_{22})\}$. Since (22) holds for each of \mathbf{s}_{i1} and \mathbf{s}_{i2} , we have undecidability (21) for both strategies of both players.

Table 5.1

	\mathbf{s}_{21}	\mathbf{s}_{22}
\mathbf{s}_{11}	$F^1(1, 1) F^2$	$(0, 1) F^2$
\mathbf{s}_{12}	$F^1(1, 0)$	$(0, 0)$

Even when G is unsolvable, there may be some case where player i has a positive decision. Table 5.1 has two subsolutions $F^1 = \{(\mathbf{s}_{11}, \mathbf{s}_{21}), (\mathbf{s}_{12}, \mathbf{s}_{21})\}$ and $F^2 = \{(\mathbf{s}_{11}, \mathbf{s}_{21}), (\mathbf{s}_{11}, \mathbf{s}_{22})\}$. Since $(\mathbf{s}_{11}, \mathbf{s}_{21})$ belongs to both subsolutions, (22) does not hold for \mathbf{s}_{i1} , but it holds for \mathbf{s}_{i2} .

Now, we show (22). A difficulty rises when the intersections of subsolutions is nonempty. Let G be any game with its subsolutions F^1, \dots, F^k . We denote the intersection $\bigcap_{l=1}^k F^l$ by \hat{F} . We stipulate that $k = 0$ and $\hat{F} = \emptyset$ if G has no Nash equilibria. If G is solvable, then $k = 1$ and F^1 is the set of all Nash equilibria $E(G)$. We note that this intersection \hat{F} satisfies interchangeability; so it can be written as $\hat{F}_1 \times \hat{F}_2$.

When all payoffs are distinct, $\bigcap_{l=1}^k F^l = \emptyset$ for $k \geq 2$. Hence, the case of $\bigcap_{l=1}^k F^l \neq \emptyset$ and $k \geq 2$ may be irrelevant from the game theoretic perspective. However, this case gives a full characterization of the existence of a positive decision, as stated in the following theorem.

Theorem 5.2. (Positive Decision) Let G be any game, $\mathbf{g} = (g_1, g_2)$ its formalized payoffs, and $i = 1, 2$. Then, for all $s_i \in S_i$, $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(I_i(s_i))$ if and only if $s_i \in \hat{F}_i$.

This has various implications: When G has no Nash equilibria, i.e., $\hat{F} = \emptyset$, $\Delta_i(\mathbf{g})$ gives no positive decisions; when G is solvable, it gives a positive decision. When G has multiple subsolutions, there are two cases; if $\hat{F} = \emptyset$, then it gives no positive decision; and if $\hat{F} \neq \emptyset$, it gives a positive decision, i.e., $s_i \in \hat{F}_i$.

The necessity in Theorem 5.2 requires a modification of the previous characterization (Theorem 4.2). We modify the target formulae $\{A_i^*(s_i)\}_{i \in S_i, i = 1, 2}$, as follows:

$$A^{**}(s_i) := \bigvee_{t_j \in \hat{F}_j} \mathbf{I}r_i^o[\text{bst}_i(s_i; t_j); \text{bst}_j(t_j; s_i)]. \quad (23)$$

This differs from $A^*(s_i)$ with the domain of disjunction \hat{F}_j instead of S_j . In this sense, it depends upon the specification of the payoff functions. We define the candidate formulae $\mathcal{C}_i = \{C_i^*(s_i)\}_{s_i \in S_i}$, $i = 1, 2$ as follows:

$$C_i^*(s_i) = \begin{cases} A_i^{**}(s_i) & \text{if } s_i \in \hat{F}_i \\ A_i^*(s_i) & \text{if } s_i \notin E(G)_i \\ I_i(s_i) & \text{otherwise.} \end{cases} \quad (24)$$

That is, $C_i^*(s_i)$ is $A_i^{**}(s_i)$ if $s_i \in \hat{F}_i$, but is $A_i^*(s_i)$ if s_i is not a Nash strategy. Crucially, it is $I_i(s_i)$ if s_i is a Nash strategy but is not a part of the intersection \hat{F}_i . The last treatment trivializes the additional premise in WF_i of (15). Then, the following characterization theorem, which will be proved in Section 5.2, implies the previous theorem and is proved before that theorem.

Theorem 5.3. (Characterization II) *Let G be any game with its subsolutions F^1, \dots, F^k , $\mathbf{g} = (g_1, g_2)$ its formalized payoffs, and $i = 1, 2$. Then, $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(I_i(s_i) \equiv C_i^*(s_i))$ for all $s_i \in S_i$.*

Theorem 5.1 also holds in $\text{EIR}^2(\mathbf{T})$, and by Lemma 2.5, $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg I_i(s_i))$ implies $\Delta_i(\mathbf{g}) \vdash \neg \mathbf{B}_i(I_i(s_i))$. Hence, in contrast to Theorem 4.4, (21) implies the incompleteness of the theory $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$. Therefore, completeness of the theory, $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$, depends upon the game G (its formalized payoffs \mathbf{g}).

Recall that, even for an unsolvable game, Theorem 3.2 states that the theory $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$ is complete within the set of game formulae. As a result, no game formulae express $I_i(s_i)$ under the theory $(\text{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$ when G is unsolvable. This observation leads us to the following theorem.

Theorem 5.4. (No-formula) *Let G be an unsolvable game, $\mathbf{g} = (g_1, g_2)$ its formalized payoffs, and $i = 1, 2$. Let $s_i \in S_i$ be a strategy for which (21) holds. Then, in $\text{EIR}^2(\mathbf{T})$, (also in EIR^2), there is no game formula A_i such that $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(I_i(s_i) \equiv A_i)$.*

5.2 Proof of the theorems

We stipulate that when $E(G) = \emptyset$, then the subsolution F is empty and $F_1 = F_2 = \emptyset$. The proof of Lemma 5.1 together with soundness for EIR^2 gives a proof of Lemma 4.2.

Lemma 5.1. *Let G be any game. Then, for any subsolution $F = F_1 \times F_2$ in G , there is a KD-model $M = (\langle W; R_1, R_2 \rangle, \tau)$ and a world $w \in W$ such that*

$$(M, w) \models \mathbf{I}_i(\mathbf{g}) \wedge \mathbf{I}_i(\mathbf{N}) \text{ and } (M, w) \models \mathbf{I}_i(\mathbf{WF}(\mathcal{A})) \text{ for all } \mathcal{A}; \quad (25)$$

$$\text{for any } s_i \in S_i, (M, w) \models \mathbf{B}_i(I_i(s_i)) \Leftrightarrow (M, w) \models I_i(s_i) \Leftrightarrow s_i \in F_i. \quad (26)$$

Proof. We construct a model $M = (\langle W; R_1, R_2 \rangle, \tau)$ satisfying (25) and (26). Let $F = F_1 \times F_2$ be a subsolution. Let $\langle W; R_1, R_2 \rangle$ be the frame given by $W = \{w\}$ and $R_k = \{(w, w)\}$ for $k = 1, 2$, i.e., it has a single world, and R_k is reflexive. Hence, this is a frame for Axiom T (and 4, 5), too. Define τ by, for $k = 1, 2$,

$$\text{for any } s; s' \in S, \tau(\text{PR}_k(s; s')) = \top \Leftrightarrow h_k(s) \geq h_k(s'); \quad (27)$$

$$\tau(w, \mathbf{I}_k(s_k)) = \top \Leftrightarrow s_k \in F_k. \quad (28)$$

That is, the preferences true relative to h_k are given by τ ; and $\mathbf{I}_k(s_k)$ is true if and only if $s_k \in F_k$. By (27), we have $(M, w) \models g_1 \wedge g_2$. Also, since $W = \{w\}$, we have, for any formula C and $k = 1, 2$,

$$(M, w) \models C \Leftrightarrow (M, w) \models \mathbf{B}_k(C). \quad (29)$$

Now, because F is a subsolution and $(M, w) \models g_1 \wedge g_2$, it follows that $(M, w) \models \text{bst}_i(s_i; s_j)$ for all $(s_i; s_j) \in F$ and for $i = 1, 2$. Thus, $(M, w) \models \text{N}0_i$. Also, $(M, w) \models \text{N}1_i$ by (28), and $(M, w) \models \text{N}2_i$ by $W = \{w\}$. Thus, $(M, w) \models \mathbf{I}_i(\mathbf{N})$ for both $i = 1, 2$.

Let us show $(M, w) \models \mathbf{I}_i(\mathbf{WF}(\mathcal{A}))$ for all \mathcal{A} . Let $\mathcal{A}_k = \{A_k(s_k)\}_{s_k \in S_k}$, $k = 1, 2$ be given. Let $E_k = \{s_k \in S_k : (M, w) \models A_k(s_k)\}$ for $k = 1, 2$. First, notice, using (29), that if $(M, w) \models \neg[\text{N}1(\mathcal{A}) \wedge \text{N}2(\mathcal{A})]$, then $(M, w) \models \mathbf{WF}_i(\mathcal{A})$. Thus, we can assume that $(M, w) \models \text{N}1(\mathcal{A}) \wedge \text{N}2(\mathcal{A})$. Using $\text{N}0_1(\mathcal{A}) \wedge \text{N}0_2(\mathcal{A})$, we have, for any $(s_1; s_2) \in S$, $(M, w) \models A_1(s_1) \wedge A_2(s_2) \supset \text{bst}_1(s_1; s_2) \wedge \text{bst}_2(s_2; s_1)$, i.e., $E_1 \times E_2 \subseteq E(G)$. Consider two cases.

(i) Let $E_1 \times E_2 \subseteq F$. Then, by (28), for $k = 1, 2$, $(M, w) \models \bigwedge_{s_k \in S_k} [A_k(s_k) \supset \mathbf{I}_k(s_k)]$; so $(M, w) \models \mathbf{WF}_i(\mathcal{A})$.

(ii) Let $E_1 \times E_2 - F \neq \emptyset$. Because F is a subsolution, it is maximal having the form of $F = F_1 \times F_2$. Also by $E_1 \times E_2 \subseteq E(G)$, we have $F - E \neq \emptyset$. Let $(s_1^*, s_2^*) \in F - E$. Then, $(M, w) \models [\mathbf{I}_1(s_1^*) \wedge \mathbf{I}_2(s_2^*)] \wedge \neg[A_1(s_1^*) \wedge A_2(s_2^*)]$ and hence for $i = 1, 2$, $(M, w) \models \neg[\mathbf{I}_i(s_i^*) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j^*) \supset A_i(s_i^*) \wedge \mathbf{B}_j(A_j^*(s_j)))]$. Thus, $(M, w) \models \mathbf{WF}_i(\mathcal{A})$ for $i = 1, 2$. ■

Proof of Theorem 5.1: Let G be an unsolvable game, and let F, F' be two subsolutions with $(s_i; s_j) \in F$ but $(s_i; s_j) \notin F'$. By Lemma 5.1, there are two models M and M' so that (25) and (26), respectively, for F and F' . Hence, $(M, w) \models \mathbf{B}_i(\mathbf{I}_i(s_i))$ but $(M', w') \not\models \mathbf{B}_i(\mathbf{I}_i(s_i))$. By soundness for EIR^2 , we have $\Delta_i(\mathbf{g}) \not\models \neg \mathbf{B}_i(\mathbf{I}_i(s_i))$ and $\Delta_i(\mathbf{g}) \not\models \mathbf{B}_i(\mathbf{I}_i(s_i))$. ■

Since the model given in Lemma 5.1 has a single world, it is a model for Axioms T, 4 and 5. Hence, Theorem 5.1 holds for EIR^2 with those axioms. In the following proof, we use the fact that Theorem 5.1 holds for $\text{EIR}^2(\text{T})$. As mentioned earlier, we first prove Theorem 5.3, followed by the proof of Theorem 5.2.

Proof of Theorem 5.3: When $s_i \in \hat{F}_i$, we have $\mathbf{I}_i^o(\mathbf{g}) \vdash A_i^{**}(s_i)$, which implies $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{I}_i(s_i) \supset A_i^{**}(s_i)$. In the other cases, by Lemma 4.1.(2), $\mathbf{I}_i^o(\mathbf{N}) \vdash \mathbf{I}_i(s_i) \supset C_i^*(s_i)$. Thus,

$$\mathbf{I}_i^o(\mathbf{g}), \mathbf{I}_i^o(\mathbf{N}) \vdash \mathbf{I}_i(s_i) \supset C_i^*(s_i) \text{ for all } s_i \in S_i. \quad (30)$$

Now, consider the converse of (30).

We modify the claims 1-3 in the proof of Theorem 4.2 as follows: for any $(s_i; s_j) \in S$,

$$(1^*): \mathbf{I}_i^o(\mathbf{g}), \mathbf{I}_i^o(\mathbf{N}) \vdash C_i^*(s_i) \wedge \mathbf{B}_j(C_j^*(s_j)) \supset \text{bst}_i(s_i; s_j).$$

$$(2^*): \mathbf{I}_i^o(\mathbf{g}), \mathbf{I}_i^o(\mathbf{N}) \vdash C_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(C_j^*(t_j)).$$

$$(3^*): \mathbf{I}_i^o(\mathbf{N}) \vdash C_i^*(s_i) \supset \mathbf{B}_j \mathbf{B}_i(C_i^*(s_i)).$$

(1*): If $C_i^*(s_i) = A_i^*(s_i)$ or $C_j^*(s_j) = A_j^*(s_j)$, then $\mathbf{I}_i^o(\mathbf{g}) \vdash \neg C_i^*(s_i)$ or $\mathbf{I}_i^o(\mathbf{g}) \vdash \mathbf{B}_j(\neg C_j^*(s_j))$; so, the assertion holds. Let $C_i^*(s_i) = A_i^{**}(s_i)$ and $C_j^*(s_j) = A_j^{**}(s_j)$. So, we have $\mathbf{I}_i^o(\mathbf{g}) \vdash \text{bst}_i(s_i; s_j)$; so, we have the assertion. Let $C_i^*(s_i) = A_i^{**}(s_i)$ and $C_j^*(s_j) = \mathbf{I}_j(s_j)$. Then, for any $k = 1, \dots, l$, $(s_i; t_j) \in F^k$ for some t_j , and also, for some k_0 , $(s_j; t_i) \in F^{k_0}$ for some t_i . Hence,

we have $(s_i; s_j) \in F^{k_0}$, i.e., $(s_i; s_j)$ is a Nash equilibrium. Hence, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \text{bst}_i(s_i; s_j)$. The case where $C_i^*(s_k) = \mathbf{I}_i(s_i)$ and $C_j^*(s_j) = A_j^{**}(s_j)$ is similar.

(2*): First, let $C_i^*(s_i) = \mathbf{I}_i(s_i)$. By N1_i, $\vdash C_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(\mathbf{I}_j(t_j))$. Then, since $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}) \vdash \mathbf{Ir}_j(\mathbf{g}) \wedge \mathbf{Ir}_j(\mathbf{N})$ by (6), we use (30) for j and get $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}) \vdash \bigvee_{t_j \in S_j} \mathbf{B}_j(\mathbf{I}_j(t_j)) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(C_j^*(t_j))$. Thus, $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}) \vdash C_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(C_j^*(t_j))$. Second, let $C_i^*(s_i) = A_i^*(s_i)$. Then, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash \neg C_i^*(s_i)$, and hence, $\mathbf{Ir}_i^o(\mathbf{g}) \vdash C_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(C_j^*(t_j))$. Third, let $C_i^*(s_i) = A_i^{**}(s_i)$. Let $s_j \in \hat{F}_j$. Then, since $\vdash \mathbf{Ir}_i^o(\text{bst}_i(s_i; s_j); \text{bst}_j(s_j; s_i)) \supset \mathbf{Ir}_j(\text{bst}_j(s_j; s_i); \text{bst}_i(s_i; s_j))$ by (6), we have $\vdash C_i^*(s_i) \supset \bigvee_{t_j \in \hat{F}_j} \mathbf{B}_j(C_j^*(t_j))$. Then, $\vdash C_i^*(s_i) \supset [\bigvee_{t_j \in \hat{F}_j} \mathbf{B}_j(C_j^*(t_j))] \vee [\bigvee_{t_j \in S_j - \hat{F}_j} \mathbf{B}_j(C_j^*(t_j))]$, equivalently, $\vdash C_i^*(s_i) \supset \bigvee_{t_j \in S_j} \mathbf{B}_j(C_j^*(t_j))$.

(3*): If $C_i^*(s_i) = A_i^*(s_i)$, we have $\vdash C_i^*(s_i) \supset \mathbf{B}_j \mathbf{B}_i(C_i^*(s_i))$ by the previous claim 3. The case for $C_i^*(s_i) = A_i^{**}(s_i)$ is similar. If $C_i^*(s_i) = \mathbf{I}_i(s_i)$, then $\vdash C_i^*(s_i) \supset \mathbf{B}_j \mathbf{B}_i(C_i^*(s_i))$ by N2_i. ■

The above three statements imply $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}) \vdash \mathbf{N}_i(C^*) \wedge \mathbf{B}_j(\mathbf{N}_j(C^*))$, and also, by (30), we have $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}) \vdash \bigwedge_{s \in S} \langle \mathbf{I}_i(s_i) \wedge \mathbf{B}_j(\mathbf{I}_j(s_j)) \supset C_i^*(s_i) \wedge \mathbf{B}_j(C_j^*(s_j)) \rangle$. Then, we using $\mathbf{Ir}_i^o(\mathbf{WF}(C^*))$, we have $\mathbf{Ir}_i^o(\mathbf{g}), \mathbf{Ir}_i^o(\mathbf{N}), \mathbf{Ir}_i^o(\mathbf{WF}(C^*)) \vdash C^*(s_i) \supset \mathbf{I}_i(s_i)$. ■

Proof of Theorem 5.2: (Only-if): Suppose $(s_i; s_j) \notin \hat{F}$ for any $s_j \in S_j$. Let s_i be not a Nash strategy. Then, $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg \mathbf{I}_i(s_i))$ by (20); so $\Delta_i(\mathbf{g}) \vdash \neg \mathbf{B}_i(\mathbf{I}_i(s_i))$ by Axiom D. Since $\Delta_i(\mathbf{g})$ is consistent by Lemma 4.2, we have $\Delta_i(\mathbf{g}) \not\vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$. Let s_i be a Nash strategy. Then, $s_i \notin F_i^l$ for some subsolution $F_1^l \times F_2^l$. Thus, $\Delta_i(\mathbf{g}) \not\vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$ by (22).

(If): If $(s_i; s_j) \in \hat{F}$ for some s_j , then $\mathbf{Ir}_i^o(\mathbf{g}) \vdash A^{**}(s_i)$. Hence, $\Delta_i^o(\mathbf{g}) \vdash \mathbf{I}_i(s_i)$ by Theorem 5.3, which implies $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$. ■

Proof of Theorem 5.4: Suppose that there is a game formula A such that $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i) \equiv A_i)$ in EIR^2 ; *a fortiori*, the same holds for $\text{EIR}^2(\text{T})$. Theorem 3.2 claims that in $\text{EIR}^2(\text{T})$, $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{B}_i(A)$ or $\mathbf{Ir}_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg A)$. This and the supposition imply $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$ or $\Delta_i(\mathbf{g}) \vdash \mathbf{B}_i(\neg \mathbf{I}_i(s_i))$ in $\text{EIR}^2(\text{T})$. This is impossible since Theorem 5.1 holds for $\text{EIR}^2(\text{T})$. ■

6 Conclusions

We have considered prediction/decision making by player i in a finite 2-person game G . We describe his decision criterion as $\mathbf{N}_i = \mathbf{N}0_i \wedge \mathbf{N}1_i \wedge \mathbf{N}2_i$ occurring in his mind, with the symmetric treatment for player j . These lead to an infinite regress of \mathbf{N}_i and \mathbf{N}_j , captured by $\mathbf{Ir}_i(\mathbf{N}_i; \mathbf{N}_j)$ in EIR^2 . We have adopted $\mathbf{Ir}_i(\mathbf{N}) = \mathbf{Ir}_i(\mathbf{N}_i; \mathbf{N}_j)$ as his basic beliefs, together with $\mathbf{Ir}_i(\mathbf{WF})$ and $\mathbf{Ir}_i(\mathbf{g})$. For a solvable game G , $\Delta_i(\mathbf{g}) = \{\mathbf{Ir}_i(\mathbf{g}), \mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF})$ determines $\mathbf{I}_i(s_i)$ as the specific formula $A^*(s_i)$ given in (12). The situation for an unsolvable G is entirely different: for some strategy s_i , $\Delta_i(\mathbf{g})$ fails to determine whether it is a possible decision or not. We discuss our game theoretic decidability and undecidability result, with comparisons to the literature as well as some possible extensions.

Positive, negative decisions, and undecidable: Suppose that G is solvable. Our game theoretic decidability result states that player i finds his Nash strategy to be a possible decision, and any non-Nash strategy to be a negative decision. Player i may find multiple possible decisions or no decisions. Our theory is silent for this choice if it exists; otherwise, negative decisions led by the emptiness may lead player i to a different decision criterion.

In contrast, when G has multiple subsolutions and hence is unsolvable, we presented the

undecidability result that player i cannot find any positive decision. One potential solution is to allow communication between the players so that they may agree upon a specific subsolution. One difficulty is that player i may not notice the necessity of this communication in the first place.

Two independent minds and discord in $\mathbf{Ir}_i(\mathbf{g})$: Theorem 5.1 is equivalent to, by Lemma 2.5, $\Delta_i(\mathbf{g}) \not\vdash \mathbf{B}_i(\mathbf{I}_i(s_i))$ and $\Delta_i(\mathbf{g}) \not\vdash \neg\mathbf{B}_i(\mathbf{I}_i(s_i))$, which is parallel to Gödel’s incompleteness theorem. Indeed, this states that the theory $(\mathbf{EIR}^2; \Delta_i(\mathbf{g}))$ (and even $(\mathbf{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$) is incomplete. These incompleteness results have some similarity but their sources are different.

Gödel’s theorem is caused by the self-referential structure of Peano Arithmetic, i.e., the theory of Peano Arithmetic can be described inside the theory itself. Our framework includes also a self-referential structure; the infinite regress operator $\mathbf{Ir}_i(\cdot; \cdot)$ includes $\mathbf{Ir}_j(\cdot; \cdot)$, and *vice versa* in \mathbf{EIR}^2 . Moreover, the criteria $\{\mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF})$ are completely symmetric between the two minds. Our undecidability arises in this context, but it is not directly generated. The direct cause lies in the infinite regress of the game $\mathbf{Ir}_i(\mathbf{g})$, which includes a possible discord between the players, depending upon whether the game is solvable or not.

Johansen’s [9] argument: He gave the following four postulates for prediction/decision making and asserted that the Nash noncooperative solution could be derived from them for solvable games.

Postulate J1 (Closed world): A player makes his decision $s_i \in S_i$ on the basis of, and only on the basis of information concerning the action possibility sets of two players S_1, S_2 and their payoff functions h_1, h_2 .

Postulate J2 (Symmetry in rationality): In choosing his own decision, a player assumes that the other is rational in the same way as he himself is rational.

Postulate J3 (Predictability): If any¹⁰ decision is a rational decision to make for an individual player, then this decision can be correctly predicted by the other player.

Postulate J4 (Optimization against “for all” predictions): Being able to predict the actions to be taken by the other player, a player’s own decision maximizes his payoff function corresponding to the predicted actions of the other player.

These postulates, except for J2, can be seen as corresponding to $\mathbf{N0}_i, \mathbf{N1}_i, \mathbf{N2}_i$ for $i = 1, 2$. Postulate J2 is interpreted as corresponding to the self-referential structure described above. That is, player i assumes the entirely symmetric structure for player j ’s thinking; Complete symmetry is obtained in terms of infinite regresses $\{\mathbf{Ir}_i(\mathbf{N})\} \cup \mathbf{Ir}_i(\mathbf{WF})$ in the logic \mathbf{EIR}^2 , while still keeping the independence of the two minds. Once $\mathbf{Ir}_i(\mathbf{g})$ is introduced, it may contain some discord. Johansen did not discuss this part.

Effective decidability of the theory: When G is a solvable game, effective decidability (decidability in the logic literature) of the theory $(\mathbf{EIR}^2(\mathbf{T}); \Delta_i(\mathbf{g}))$ follows from the full completeness theorem (Theorem 4.2). For $(\mathbf{EIR}^2; \Delta_i(\mathbf{g}))$, we need to restrict the class of formulae. When G is unsolvable, this argument does not work: the effective decidability in such a case remains open.

Other variants and extensions: Our results (3) and (4) are obtained for \mathbf{EIR}^2 . As stated, those results hold for a stronger system than \mathbf{EIR}^2 , for example, in those with any of Axioms T, 4, and 5, but we choose \mathbf{KD}^2 to keep subjectivity of each player. In the present logic \mathbf{EIR}^2 , player

¹⁰This “any” was “some” in Johansen’s original Postulate 3. He assumed (p.435) that the game has the unique Nash equilibrium. In this case, the above difference does not matter.

i has the theory $\Delta_i(\mathbf{g})$, but player j can have his own theory $\mathbf{B}_j(\Gamma_j)$, which may be entirely different from $\Delta_i(\mathbf{g})$. If they recommend compatible decisions and predictions, the players may not find the differences in their theories by watching the *ex post* play. This is not allowed in the logic $\text{EIR}^2(\text{T})$. Thus, EIR^2 enables us to separate between subjective thinking and actual plays. This separation may deserve further investigation.

We have confined ourselves to the 2-person case both for the logic and game theory. For n -person case ($n \geq 3$), we would meet new problems in both epistemic logic and game theory. We will discuss those extensions in separate papers.

Other game theoretic undecidability: Kaneko-Nagashima [10] gave a 3-person game having a unique Nash equilibrium in mixed strategies. It is assumed that the game structure and real number theory Φ_{rcf} (real closed field theory) are common knowledge among the players in an infinitary predicate logic. They showed that $\mathbf{C}(\exists x \text{Nash}(x))$ is provable from their common knowledge of G and Φ_{rcf} , but that neither $\exists x \mathbf{C}(\text{Nash}(x))$ nor $\neg \exists x \mathbf{C}(\text{Nash}(x))$ is provable. That is, the players commonly know the abstract existence of a Nash equilibrium, but do not find a concrete one; hence they cannot play the specific Nash equilibrium strategy.

This undecidability is caused by the lack of names for some irrational numbers such as $\sqrt{51}$ in their language, which is involved in the Nash equilibrium in the 3-person game with rational payoffs. The main reason for this difficulty is to give a name to a concept, but not the self-referential structure.

Other game theoretic solution concepts: The game theory literature has various “solution concepts” other than the Nash theory. One concept is the “dominant strategy” criterion, which requires a player to choose one which is best against any strategy of the player. We can extend this by requiring one player to use a best response against any dominant strategy of the other player, predicting that the other player adopts the dominant strategy criterion. Even we can extend this argument to any finite level. In those cases, we have game theoretic decidability result. We conjecture that any solution concept which does not require infinite regress will lead to similar decidability.

Future directions: Our approach assumes unbounded logical abilities and unbounded interpersonal thinking, but we still meet the undecidability result. From the social science perspective, it may be fruitful to investigate whether a theory with bounded logical abilities or bounded interpersonal thinking can avoid undecidability.

References

- [1] Blackburn, P., M. de Rijke, and T. Venema, (2001), *Modal Logic*, Cambridge University Press, Cambridge.
- [2] Boolos, G., (1979), *The Unprovability of Consistency*, Cambridge University Press, Cambridge.
- [3] Brandenburger, A., (2014), *The Language of Game Theory*, World Scientific, London.
- [4] Fagin, R., J. Y. Halpern, Y. Moses and M. Y. Verdi, (1995), *Reasoning about Knowledge*, The MIT Press, Cambridge.

- [5] Heifetz, A., (1999), Iterative and Fixed Point Common Belief, *Journal of Philosophical Logic* 28, 61-79.
- [6] Hu, T., and M. Kaneko (2012), Critical Comparisons between the Nash Noncooperative Theory and Rationalizability, *Logic and Interactive Rationality Yearbook 2012*, Vol.II, eds. Z. Christo, *et al.* 203-226, http://www.ilic.uva.nl/dg/?page_id=78
- [7] Hu, T., and M. Kaneko (2014), Epistemic Infinite Regress Logic, to be completed in 2014.
- [8] Hu, T., M. Kaneko, and N.-Y. Suzuki, (2014), Small Infinitary Epistemic Logics and Some Fixed-Point Logics, to be completed in 2014.
- [9] Johansen, L., (1982), On the Status of the Nash Type of Noncooperative Equilibrium in Economic Theory, *Scand. J. of Economics* 84, 421-441.
- [10] Kaneko, M., and T. Nagashima, (1996), Game logic and its applications I, *Studia Logica* 57, 325-354.
- [11] Kaneko, M., (2002), Epistemic logics and their game theoretical applications: Introduction. *Economic Theory* 19, 7-62.
- [12] Kaneko, M., (1999), Epistemic considerations of decision making in games. *Mathematical Social Sciences* 38, 105-137.
- [13] Kline, J. J., (2013), Evaluations of epistemic components for resolving the muddy children puzzle, *Economic Theory* 53, 61-84.
- [14] Meyer, J.-J. Ch., van der Hoek, W., (1995), *Epistemic logic for AI and computer science*. Cambridge.
- [15] Mendelson, E., (1988), *Introduction to Mathematical Logic*, Wadsworth, Monterey.
- [16] Nash, J. F., (1951), Non-cooperative Games, *Annals of Mathematics* 54, 286-295.
- [17] Osborne, M., and A. Rubinstein, (1994), *A Course in Game Theory*, MIT Press, Cambridge.
- [18] Suzuki, N.-Y., (2013), Semantics for intuitionistic epistemic logics of shallow depths for game theory, *Economic Theory* 53, 85-110.
- [19] Van Benthem, *Logic in Games*, Institute for Logic, Language and Computation.
- [20] Van Benthem, J., E. Pacuit, and O. Roy, (2011), Toward a Theory of a Play: A Logical Perspective on Games and Interaction, *Games* 2, 52-86.