# Logical Monotony and Bounded Rationality

## Mikael Cozic, Paris IV-Sorbonne

3 janvier 2004

The starting point of my actual research is the problem of logical omniscience and other similar idealizations that one can find in epistemic and doxastic models like epistemic logic, probability theory or belief revision; these idealizations are expressed by the following inference rules :

epistemic logic	$\frac{A \to B}{K_i A \to K_i B}$	(REM)
probability	$\frac{A \to B}{L_{i,\alpha} A \to L_{i,\alpha} B}$	(RPM)
belief revision	if $\vdash_{PC} A \to B$ , then if $A \in K_i, B \in K_i$	(PMR)

I call the existence of such idealizations the phenomenon of logical monotony. Lots of work has been done to weaken logical monotony in epistemic logic; one can find a good survey in Reasoning about Knowledge by Fagin and *al*.

Some work has been done recently to weaken logical monotony in belief revision, for example by R.Wasserman and R.Parikh.

One disputed question is whether or not logical monotony is an innocuous idealization.

If one adopts a realist view of epistemic models, it is clearly not.

But if one adopts a more instrumentalist view, the answer to this question depends on the role of epistemic models in more general settings like formal semantics, distributed systems or rational choice theory; for short, it depends on the application's field.

I won't say anything concerning the first two fields, but I will argue that the weakening of logical monotony is relevant in the case of rational choice theory. I call this methodological thesis the pragmatic relevance of logical monotony.

Basically, it seems to me that the pragmatic relevance of logical monotony comes from two essential features of rational choice theory :

- on the one hand, rational choice theory makes an intensive use of epistemic models and inherits their logical monotony

- on the other hand, every rational behavior can be seen as the result of a (not necessarily conscious) reasoning, and more precisely an optimizing reasoning

The consequence of those two features is that at the same time practical reasoning is a central component of decision making but its complexity is omitted.

In few words, rational choice theory can take into account only information ascribed to agents, and not the deductive processing of this information.

Let's be a little bit more precise now.

In a decision making, there are different levels at which deductive reasoning can be relevant :

(1) at the level of representation, an agent can have a partial representation of the environment because of a lack of deductive processing :

(1a) it can be the case that the knowledge of one feature of the environment depends on the ability to solve an explicit mathematical question

for a first example, taken from [Sa 67], let's think to a bet on a mathematical question, like : "Is d the n-th digit of  $\pi$ ?"

for a second and more realistic example, it is sufficient to consider a decision problem that involves public-key cryptography like the Rivest-Shamir-Adleman system.

Suppose that

- the agent is the executive council of the UN about the Irakian crisis

- the opportunities of the council are : to accept the attack of Irak by USA or : to condemn the attack

- the environment is reduced to an encoded message from a person who knows exactly the military resources of Irakian Army and the (public) encoding key of the receiver of the message

If the council wants to know the content of the message, the decoding key must be guessed; but this key can be guessed from the encoding key only if one is able to factor efficiently a very large integer; today, no one can do this; so there is one (very important!) feature of the environment that is ignored because of a lack of deductive ability.

(1b)it can be the case that the knowledge of one feature of the environment depends on the ability to deduce any consequence from the available informations

let's consider for example the game of chess; since Zermelo 1913, one knows that,

provided that if a position is repeated three times, the outcome is a draw, the game of chess is finite and has a value that is

- either there exists a strategy for white that ensures that white wins

- or there exists a strategy for black that ensures that black wins

- or every player has a strategy that ensures a draw

Then, under the logical monotony assumption, as H.Simon has noted "chess is a trivial game"([Si 92])

(2) at the level of optimization, an agent can be ignorant of the optimal action, given his preferences and his information

It is worth noting that the environment can be relatively simple (= part of the data of the optimization problem) and the choice of the optimal action (solving the optimization problem) very hard :

for example, in a combinatorial auction, the auctioneer has to select an optimal allocation given the bids of the buyers; but this problem is computationaly as hard as the well-known MAX-CLIQUE problem, which is NP-Hard and even hard to approximate. This is intensively investigated today by the so-called "computational mechanism design" pioneered by Noam Nisan and others.

(3) at the level of execution, as soon as the strategy is a decision function which depends on the past action of others and/or of nature, the construction of such a function and the computation following such a function can be problematic

example : in AI, real-time tasks where the computation following the optimal decision function can violate the deadline

Even if these distinctions between the component of representation and the component of optimization among logical difficulties of a decision problem has no sharp boundaries, it is clear that those difficulties stresses an important difference : the difference between

- "available" informations is an information that effectively helps the agent to decide

- "virtual" informations is informations that the agent could have if he would process perfectly all that he knows

cf [Sta99] For example, the value of chess and the knowledge of irakian military resources in the previous scenario are only virtual informations.

When agents have to act on the basis of their informations, the relevant notion of information is that of available information; and as far as by solving the problem of logical monotony we are trying to better capture this availability of information, the problem of logical monotony is pragmatically relevant.

The connexion with bounded rationality is quite evident now : bounded rationality, in the general sense, is the research program whose task is to refine basic assumptions of the rational choice model. Since H.Simon, one of the main direction was the revision of cognitive and especially computational assumptions; according to him,

"Broadly stated, the task is to replace the global rationality of economic man with a kind of behavior that is compatible with the access to information and the computational capacities that are actually possessed by organisms, including man..." ([Si55]) The weakening of logical monotony is closely related to this "computational part" of the Bounded Rationality program, since its purpose is to take into account our deductive limitations.

But once the pragmatic relevance of logical monotony is accepted, the main question is :

what could and what should we expect from a weakening of logical standards?

= how could the weakening of logical standards help to better reflect information's processing and its availability for action?

I won't give any general answer to this question, but in the second part of my talk, I would like to give some insights into different ways of fulfilling this program by considering game-theoretic examples.

I've chosen game-theoretic examples because, as noted by Robert Stalnaker, "the game situation is a central case where the problem of logical omniscience is particularly acute" ([St98]).

When one wants to entertain such a program, typical basic methodological directions depends on questions like the following ones :

- is it better to make general weakenings (for example, restrict some functions of the decision making to computable functions) or more context-dependant weakenings (depending on the type of games, or on the type of environment, etc...)?

- is it better to build a measure of weakening (for example to take a complexitytheoretic measure) or simply to drop or replace the deductive assumptions?

- should we look for the emergence of "positive" results from the weakenings, or can they only refute preexisting results that were based on strong cognitive assumptions?

For example, I've recently worked on the well-known computational restrictions on the strategies in the Repeated Prisoner's Dilemma; those restrictions are weakenings that

- depend closely on the context of repeated games

- build a measure based on the size of the minimum automata

- permit to reach (approximately) the cooperation's payoff, what is impossible either in the one-shot Dilemma, or in the Finitely Repeated Dilemma

actually, the basic result is that if one of the strategies is subexponentially restricted, the cooperation's payoff can be reached

Today, I will confine myself to show by an example that one can already do quite interesting things even with a very general method of weakening.

Among the "epistemic models", epistemic logic is probably the most simple and the one in which solutions to the logical monotony problem are the most investigated. The main use of epistemic logic in game theory concerns the cognitivebayesian justification of solution concepts like Nash Equilibrium, Correlated Equilibrium, etc...The basic motivation of this field is that whereas it is rational for a player to stay in a Nash Equilibrium once it is reached, it is much less evident to understand why rational players should play an action profile that is a Nash Equilibrium when they ignore the choices of the others.

The cognitive justification of a solution concept Eq is a set of doxastic conditions such that if all are satisfied by bayesian players, these players choose an action profile in the solution concept

Let us look at a solution concept that is quite economical in doxastic conditions, namely the Rationalizability; it corresponds to the iterated elimination of strongly dominated strategies.

For example, let us start from this two-player game ([BatB]);

1 / 2	а	b	с
A	(3,0)	(1,0)	(0,1)
В	(1,1)	(0, 2)	(1,1)
С	(0, 0)	(4, 1)	(2,2)
D	(0,3)	(1,0)	(3,2)

B is s.dominated by (1/2 A, 1/2 D)

1 / 2	a	b	с
А	(3,0)	(1,0)	(0,1)
С	(0,0)	(4,1)	(2,2)
D	(0,3)	(1,0)	(3,2)

1/2	a	c	1 / 0	1	<u> </u>
Å	(3.0)	(0.1)	1/2	a	0
Λ	(0,0)	(0,1)	A	(3.0)	(0.
C	(0,0)	(2,2)		(0,0)	$(\circ,$
D	(0.3)	(3.2)		(0,3)	(3,
	(0,3)	(0,2)			

At the first step, one eliminates strategies such that no conjecture on the strategies of the other player could make it a best response. Here the strategy B of player 1 is strongly dominated by the mixed strategy (1/2 A, 1/2 D) Then, in the new game, one iterates the elimination, and so on...

It is important to stress that this solution concept is intended to capture the "*logical consequence*" ([BD 87]) of common knowledge of the rationality of the players and of the structure of the game (set of players, possible actions, and preferences).

To give some doxastic conditions to the equilibrium, the canonical method associates an epistemic model to games; those models can be "type spaces" - following the technique of Harsanyi ([TT 88], [BatB 99]) - or "belief system"([AH 02], [St 94], [BatB 99]).

A belief system is just a set of states with, for each player, in each state, a probability distribution on the set of states.

★ a **belief system** for a set N of agents is a n-uple  $\mathbf{B} = \langle N, S, (p_i)_{i \in N} \rangle$ where

 $\bullet~S$  is a finite set of states

• for every player  $i, p_i : S \to \Delta(S)$  associates to every state a probability distribution on S s.t.

if 
$$p_i(s, \{s'\}) > 0$$
, then  $p_i(s) = p_i(s')$  (ND)

One can remark that a belief system induces an epistemic logic's structure when one defines the accessibility relation as the support of the probability distribution.

\* for every player  $i \in N$ , for every state  $s \in S$  the probability distribution  $p_i(s)$  induces an **accessibility relation**  $R_i$  defined as below :

$$sR_is'$$
 iff  $s' \in Supp(p_i(s))$ 

Obviously, the structure induced by a belief system is a KD45 structure.

Let G be a game :

 $\star$  a finite normal-form game for N-players is a n-uple  $G = < N, A = (A_i)_{i \in N}, u = (u_i)_{i \in N} >$  where

- $\bullet~N$  is the set of players
- $A_i$  is the set of actions of player i
- $u_i: A \to \mathcal{R}$  is the pay-off function of player *i*

A belief system is associated to a game G when an action profile is assigned to each state of the belief system :

★ a belief system **B** is **associated** to the game *G* if there exists a collection of functions  $(\sigma_i)_{i \in N}$  where  $\sigma_i : S \to A_i$ 

In each state s of a belief system, each player chooses an action and has an induced conjecture ie a probability distribution on the strategies of other players. One supposes that each player maximizes his expected utility given his conjecture. The search for cognitive justification consists in finding to which (doxastic) conditions the action profile chosen in s is in the solution concept.

To do this, one can build a very simple formal language with atomic formulas

- P which means morally "the profile is rationalizable"
- $RAT_i$  which means morally "player *i* is rational"

and then defines the notion of a **normal** interpretation which captures the intended meaning of the atomic formulas :

(1) 
$$s \models P$$
 iff  $(\sigma_1(s), ..., \sigma_n(s)) \in A^{\infty}$  where  

$$-A_i^0 = A_i \text{ et } A^0 = \prod_{i \in N} A_i^0$$

$$-A_i^{k+1} = \{a_i \in A_i^k : \exists \mu_i \in \Delta(A_{-i}, a_i \in r_i(\mu_i), \mu_i(A_{-i}^k) = 1\}$$

$$-A_i^{\infty} = \bigcap_{k \in N} A_i^k \text{ et } A^{\infty} = \prod_{i \in N} A_i^{\infty}$$
(2)  $s \models RAT_i$  iff for every  $a_i \in A_i$ ,  
 $\sum_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i}) mrg_{A_{-i}} p_i(s, [a_{-i}]) \leq \sum_{a_{-i} \in A_{-i}} u_i(\sigma_i(s), a_{-i}) mrg_{A_{-i}} p_i(s, [a_{-i}])$ 

[!] that is s satisfies the formula  $RAT_i$  iff the action taken by player i in s is optimal given his marginal probability on the actions of other players

(3) propositional constants have their usual set-theoretic definitions

(4)  $s \models B_i A$  iff for every s' s.t.  $sR_i s', s' \models A$ 

(5)  $s \models CA$  iff for every s' in the universal closure of  $R_i, s' \models A$ 

Here comes the basic result :

**Theorem** ([TW 88], Theorem 5.2; [St 94], Theorem 1; [BatB 99], Proposition 3.10)

Let a game G, B a belief system associated to G and I interpretation normal in each state of B; then in every  $s \in S$ ,

if for every  $i, s \models (RAT_i \land CRAT_i)$ , then  $s \models P$ 

**Question** What happens now if we do not make any more the assumption that players are perfect reasoners?

Let's take one of the weakened epistemic logic, the semantics of "impossible possible worlds" and let's see if the result still holds.

Of course, we're starting from a true possible (possible for brief) state, ie a state where holds a normal interpretation,

let a game G, a belief system **B** associated with G and I an interpretation; a state  $s \in S$  is **possible**w.r.t. I if I is normal in s

and we're using impossible world only in the accessibility relations of the agent.

It then trivial to show that in a possible state, there can be rationality and common knowledge of rationality without the actual play of an action profile that survives the iteration of strongly dominated actions.

Fact 1

there exists a game G, a belief system **B**, an interpretation I and a *possible* state  $s \in S$  where it is not true that if for every  $i, s \models (RAT_i \land CRAT_i)$ , then  $s \models P$ 

Proof

Let's consider the previous game et let the following belief system :  $S = \{Cb, Cc, Ac\}$   $p_1(Cb, \{Cc\}) = p_1(Cb, \{Cb\}) = p_1(Cc, \{Cc\}) = p_1(Cc, \{Cb\}) = 1/2$   $p_2(Cc, \{Cc\}) = p_2(Cc, \{Ac\}) = p_2(Ac, \{Cc\}) = p_2(Ac, \{Ac\}) = 1/2$  $\sigma(Cc) = (C, c), \ \sigma(Ac) = (A, c), \ \sigma(Cb) = (C, b)$ 

Let's define the interpretation  $I^* : I^*$  is normal in Cc but we stipulate that for every i;  $Ac, Cb \models RAT_i$ . One can then check that  $I^*$  is not normal in Cband Ac and that

 $Cc \nvDash P$ 

 $Cc \models RAT_1$  (the expected utility of player 1 is 3)

 $Cc \models RAT_2$  (the expected utility of player 2 is 3/2)

 $Cc \models CRAT_i$  for i = 1, 2

#### Remark

One can ask : how is this kind of example exactly connected with (weakening of) logical monotony?

1.

Let's enrich the language with two kinds of atomic propositions : - BEST(a, i) means that action a is a best response of player i- PLAY(a, i) means that action a is chosen by player i

Clearly, in every belief system, in every state  $\boldsymbol{s}$  and for every normal interpretation  $\boldsymbol{I}$ 

$$s, I \models (-BEST(a, i) \land PLAY(a, i)) \rightarrow -RAT_i (R)$$

Let's define  $\models \varphi \rightarrow \psi$  the fact that in every state s of every belief system, if I is normal,  $s, I, \models \varphi \rightarrow \psi$ .

In this context, logical monotony means that in every state s and for every normal interpretation I, if  $s, I \models B\varphi$  and if  $\models \varphi \rightarrow \psi$ , then  $s, I \models B\psi$ .

If one extends  $I^*$  in the (normal) following way :

 $Cb \models -BEST(b, 2)$  $Cb \models PLAY(b, 2)$ 

then the lack of omniscience of player 1 is made explicit : if  $\varphi := (-BEST_{(b,2)} \land PLAY_{(b,2)})$  and  $\psi := -RAT_2$ , then logical monotony would imply that

$$s, I^* \models B_1(-\psi) \to B_1(-\varphi)$$
, hence  
 $s, I^* \models -B_1(-\varphi) \to -B_1(-\psi)$   
But here we have :  
 $s, I^* \models -B_1(-\varphi) \to B_1(-\psi)$ 

s,  $I^* \models -B_1 - (-BEST_(b, 2) \land PLAY_(b, 2))$  and  $s, I^* \models B_1(RAT_2)$ 

2.

This failure of omniscience can be connected with the algorithm of iterated elimination of strongly dominated strategies : at step 1, B is eliminated : then, at step 2, b is eliminated. But in the previous model, b is considered by player 1 as a possible action for player 2 in Cc. Therefore, one can interpret the failure of logical omniscience in the model (as showed above) as the failure to do step 2 in the algorithm.

#### Interpretation of Fact 1

This result has very few no technical interest. It means before all that with a simple generalization of epistemic logic, one can model the pragmatic effect of a limited ability to draw logical consequences of one's beliefs.

But there is something more interesting : from this example, one can argue that the impossible possible worlds semantics is especially well suited to model this kind of situation. Why? Because there is a strong correspondence between epistemic elimination of impossible states and pragmatic elimination of dominated actions, ie a strong correspondence between

- the fact that in the belief system some agents do not eliminate (in their accessibility relation, then "epistemically" eliminate)) impossible states

- the fact that some strategies should be eliminated (in the sense of "iterated elimination") if the agents were able to draw all the consequences of their knowledge, but in fact are not eliminated

Now, it is a general feature of most of models of rational choice that the choice of an agent is represented as the elimination of all possible actions, except one; so one has good reasons to think that this correspondence between doxastic and pragmatic elimination make the impossible possible world well suited in general.

In fact, if we go back to the initial example, one can do more with the simple tool of impossible possible worlds than only refute the classical justification of Rationalizability. One can indeed

- build a measure of the global logical ability of the players

- regain in a controlled way the classic justification on the basis of the mentioned measure

How to do this?

First, I will slightly enrich the language with epistemic operator of truthful iterated mutual belief and with atomic propositions to express the different steps of iterated elimination :

\* let  $L_2$  multi-epistemic propositional language enriched by a common knowledge operator C by truthful iterated mutual belief operator  $\{M^k : k \in N\}$  and based on the set of atomic propositions  $At_2 = At_1 \bigcup \{P^h : h \in N\}$ 

 $\star$  the extension of the notion of normal interpretation to  $L_2$  is straightforward :

(6) for every  $\varphi \in Prop(L_2)$ ,

$$\begin{split} s &\models M^0 \varphi \text{ iff } s \models \varphi, \\ s &\models M^1 \varphi \text{ iff } s \models \varphi \land (B_1 A \land B_2 \varphi), \\ s &\models M^k \varphi \text{ iff } s \models (\bigwedge_{l=0}^{l=k-1} M^l \varphi) \land (B_1 M^{k-1} \varphi \land B_2 M^{k-1} \varphi) \end{split}$$

that is at s,  $\varphi$  is true and each player believes that it is true and believes that the other believes that it is true and so on until level k.

(7)  $s \models P^h$  iff  $\sigma(s) \in A^h$ 

that is the profile played in s survives h eliminations

I have previously defined the fact that a state is possible if the interpretation in it is normal; now, one can define for a state s the fact that he is k-normal, ie the fact that all states accessible in 0 to k steps from s is normal. This will be the measure of logical ability.

\* let s a state of a belief system associated with a game G; an interpretation is (s, k)-normal,  $k \in \mathbb{N}$ , if for every state accessible from s in 0 to k steps by the accessibility relations of players is possible.

From this notion, one can deduce a (still trivial) Fact : if a state is k-normal and if agents in the state have mutual belief of level k that they are rational, then the actions played at state s survive iterated elimination until level k + 1. Fact 2

Let  $s \in S$  a state of a belief system associated with G and I a (s, k)-normal interpretation; if for every  $i, s \models M^k RAT_i$ , then  $s \models P^{k+1}$ 

Final remarks :

and

(1) Since the set of actions of each agent is finite, there is a level h s.t. the set of actions that survive the iteration of level h is the same that the set of actions that survive all the iteration. Let  $h^*$  the lowest such level.

So, if  $k = h^* - 1$ , then the action profile played will be rationalizable. Consequently, if k gives the measure of the global logical ability of the agents,  $k-h^*-1$  gives the distance between the logical ability of agents and the minimum logical ability sufficient to reach a rationalizable profile.

(2) if the role of epistemic logic is

"to determine what assumptions on the beliefs and reasoning of the players are **implicit** in various solutions concepts" (Bonanno, "Modal Logic and Game Theory"),

one can see the use of impossible states as a slight improvement because now the assumptions on reasoning's abilities are explicit.

(3) the improvement comes from the correspondence between

- (epistemic) non-elimination of impossible states

- (pragmatic) non-elimination of non-rationalizable strategies

(4) but as such, impossible states doesn't give any predictive power; the best they can do is to give a conditional predictive power : if one has reason to think that he is modelling agents whose effective reasoning corresponds to level k, then he will only eliminate the strategies eliminated by the first k + 1 iteration.

### Références

[AB 95] R.Aumann et A.Brandenburger, "Epistemic Conditions for Nash Equilibrium", *Econometrica*, vol.63, 5, pp. 1161-1180

[AH 02] R.Aumann et A.Heifetz, "Incomplete Information", dans R.Aumann et S.Hart, Handbook of Game Theory, vol.III

[Ber 84] D.Bernheim, "Rationalizable strategic behavior", *Econometrica*, vol.52, 4, pp. 1007-1028

[Bra 92] A.Brandenburger, "Knowledge and Equilibrium in Games", *The Journal of Economic Perspectives*, vol.6, 4, pp. 83-101

[BraD 87] A.Brandenburger et E.Dekel, "Rationalizability and Correlated Equilibria", *Econometrica*, vol.55, 6, pp. 1391-1402

[BatB 99] P.Batigalli et G.Bonanno, "Recents results on belief, knowledge and the foundations of game theory", *Research in economics*, 53, pp. 149-225

[P 84] D.Pearce, "Rationalizable strategic behavior and the problem of perfection", *Econometrica*, vol.52, 4, pp. 1029-1050

[St 94] R.Stalnaker,"On the Evaluation of Solution Concepts, *Theory and Decision*, vol.37, 1, pp. 49-73

[St 99] R.Stalnaker, "The Problem of Logical Omniscience : II", in <u>Context and Content</u>, OUP, 1999

[TW 88] T.Tan et S.Werlang, "The Bayesian Foundations of Solution Concepts of Games", *Journal of Economic Theory*, vol.45, 2, pp. 370-391

 $[W\ 02]$  B. Walliser, "Les justifications des notions d'équilibres de jeux", à paraître dans  $Revue\ d'Economie\ Politique$