

An Epistemic Generalization of Rationalizability

Extended abstract

Rohit Parikh

*City University of New York
April 14, 2016*

Abstract: *Savage showed us how to infer an agent's subjective probabilities and utilities from the bets which the agent accepts or rejects. But in a game theoretic situation an agent's beliefs are not just about the world but also about the probable actions of other agents which will depend on their beliefs and utilities. Moreover, it is unlikely that agents know the precise subjective probabilities or cardinal utilities of other agents. An agent is more likely to know something about the preferences of other agents and something about their beliefs. In view of this, the agent is unlikely to have a precise best action which we can predict, but is more likely to have a set of "not so good" actions which the agent will not perform.*

Ann may know that Bob prefers chocolate to vanilla to strawberry. She is unlikely to know whether Bob will prefer vanilla ice cream or a 50-50 chance of chocolate and strawberry. So Ann's actions and her beliefs need to be understood in the presence of such partial ignorance. We propose a theory which will let us decide when Ann is being irrational, based on our partial knowledge of her beliefs and preferences, and assuming that Ann is rational, how to infer her beliefs and preferences from her actions.

Our principal tool is a generalization of rational behavior in the context of ordinal utilities and partial knowledge of the game which the agents are playing.

1 Introduction

Suppose that a man is hungry (desire), knows that there is a restaurant within two blocks (knowledge) and is able to walk the two blocks (ability). Then we can predict that he will go to that restaurant. Knowledge plus desire plus ability lead to action.

Suppose that there are two restaurants nearby and we do not know his taste. Then we can predict that he will go to one, but not know which one.¹ But we know that he will go to one of them. The strategy of not going to either is dominated for him by going to restaurant A as well as by going to restaurant B.² We can rule out the dominated strategy of not going anywhere. (I am assuming that he knows also that the restaurants are open and affordable to him).

Suppose now that *we* know that the man is hungry, can walk two blocks and also know that he does not know about the restaurant. Then we will *tell* him about the restaurant. If we are nice people then the strategy of not telling him anything is dominated by telling him about the restaurants. If we do not tell him and he finds out later that we knew then he will reproach us. So telling is better.

If on the other hand we know he already knows about the restaurant then it would be rude to tell him anyway. He might be offended.

Thus our actions take place in a world of desires, knowledge (or beliefs) and abilities. And quite often not only our own beliefs and desires are involved but also what we know about the desires and beliefs of others.

Leonard Savage [16] worked out a theory in which by observing an agent's willingness to accept or reject certain bets we can discover both his beliefs (subjective probability) and his desire (utility). This theory has been questioned and has some difficulties pointed out by Ellsberg, Allais and Kahneman and Tversky. But the theory is much respected and still taught routinely.

But Savage did not have a theory of what we do when other agents are involved and we know *something but not everything* about their desires and

¹It may be that he prefers restaurant A in which case the strategy of going to B is dominated. But *we* do not know that.

²Intuitively a strategy *s* is dominated by *s'* if the outcome from *s'* is always better than the outcome from *s*. And here we interpret "always" as "in all situations compatible with our knowledge." Thus the notion of *dominated* has an epistemic (or doxastic) component.

beliefs. Moreover, we are more likely to know ordinal utilities than cardinal utilities for other agents.

Two tools have come into prominence since Savage. One is *Game Theory* which deals with many agents taking actions, taking into account the other agents' desires. Typically the situation is described either in terms of a payoff matrix for a normal form game or in terms of a tree of actions and choices in the case of an extensive form game.

Another tool is *Epistemic Logic* which allows us to represent the states of knowledge or belief of many agents. We can represent "Ravi knows that it is raining" or "Yang does not know that Ravi knows that it is raining" using Kripke structures as well as Aumann structures [1].

However, these two tools have not been brought together in the way we need. Of course there is the field of Epistemic Game Theory but it investigates certain models originally developed by Harsanyi and investigated further by Brandenburger and others. See [3] for example. The more *down to earth* area which we are proposing does not exist.

Here is a story about the TV detective Adrian Monk. A woman has fallen off a fifteenth floor balcony and is lying on the pavement, dead of course. A policeman arrives and a little later Adrian Monk arrives. "We don't know yet if it was murder or suicide" says the policeman.

"It was murder," says Monk.

"But you just arrived! How can you know?"

"She was in the middle of painting her nails," says Monk.

Now we all agree that a woman does not paint half her nails and then commit suicide without painting the other half. Perhaps she was in despair because she ran out of paint but that is not very likely. Murder is the far more plausible explanation. Since Monk is unable to answer the question, "Why might she commit suicide at *that* moment?" he concludes that it was murder.

When people or animals or children do something, we look for an explanation in terms of their beliefs and desires. Often the desires are known to us. A hungry animal wants to eat. A tired child wants to go to bed. And then from their actions we conclude what their beliefs are. Sometimes we go in the other direction and manipulate their beliefs so as to get them to act

in a certain way. “If you are a good girl, you will get ice cream after the dinner!”

This paper is about inferring beliefs from actions, and also, inferring preferences from actions. We offer some examples from literature, from real life and offer a formal framework. But we will be guided by the following intuition. If, given an agent’s desires and beliefs, action s is definitely worse than action s' then the agent will not do action s . If the agent *does* do action s then we can conclude that we were wrong about the beliefs or desires.

Of course the agent might be irrational. All of us have met irrational people. But the scenarios we consider are so simple that irrationality is unlikely to be an explanation.

With language using creatures, we do sometimes just ask, to find out their beliefs. But this method is not available with small children or animals. And it may not work even with adults who might have a reason to deceive. So a formal theory of inferring beliefs from actions is likely to have value.

2 Inducing Beliefs

2.1 Shakespeare’s *Much ado about Nothing*

At Messina, a messenger brings news that Don Pedro, a Spanish prince from Aragon, and his officers, Claudio and Benedick, have returned from a successful battle. Leonato, the governor of Messina, welcomes the messenger and announces that Don Pedro and his men will stay for a month.

Beatrice, Leonato’s niece, asks the messenger about Benedick, and makes sarcastic remarks about his ineptitude as a soldier. Leonato explains that “There is a kind of merry war betwixt Signior Benedick and her.”

Various events take place and Claudio wins the hand in marriage of Hero, Leonato’s only daughter and the wedding is to take place in a week.

Don Pedro and his men, bored at the prospect of waiting a week for the wedding, hatch a plan to matchmake between Beatrice and Benedick who inwardly love each other but outwardly display contempt for each other.

According to this strategem, the men led by Don Pedro proclaim Beatrice’s love for Benedick while knowing he is eavesdropping on their conversation. Thus we have, using K for “knows”, b for Benedick, d for Don Pedro and E for the event of eavesdropping,

$$K_b(E), K_d(E) \text{ and } \neg K_b(K_d(E))$$

Both Benedick and Don Pedro know about the eavesdropping but Benedick does not know that Don Pedro knows.³

All these conditions are essential and of course the plot would be spoiled if we had $K_b(K_d(E))$ instead of $\neg K_b(K_d(E))$. Benedick would be suspicious and would not credit the conversation.⁴

The women led by Hero carry on a similar charade for Beatrice.

Beatrice and Benedick are now convinced that their own love is returned, and hence decide to requite the love of the other.

The play ends with all four lovers getting married.

Benedick's Decision problem

	love	nolove
propose	100	-20
nopropose	-20	0

Here *love* means “Beatrice loves me” and *nolove* the other possibility.

The payoffs are explained by the fact that if he proposes and she does not love him, she will ridicule him, so -20. If she does love him then proposing has a big payoff of 100 and not proposing has a payoff of -20 for opportunity lost.

As long as he believes that Beatrice does not love him, not proposing dominates proposing in terms of payoffs.

Once he realizes that she does love him, proposing becomes dominant.

Note that proposing always *was* his dominant strategy. She did love him. But his false belief guided his actions.

Was he rational when he did not propose to her even though given the situation, proposing was his dominant strategy? We would have to say that he was rational and that we need to define his rationality *relative to his knowledge*.

³Here K stands for knowledge, the subscript b denotes Benedick and E denotes the event of eavesdropping. \neg stands for negation.

⁴If $K_b K_d(E)$ then Benedick would say to himself, “If he knows I am here why not talk to me directly? What devious motive can he have?”

2.2 Benedick and the two florists - going to second order

Suppose there are two florists in Messina. If there is a wedding they will have to furnish flowers and that many flowers are only available in Napoli. So to get flowers for a wedding they have to write to Napoli and put down a deposit.

They both know that Beatrice loves Benedick. But they do not both know whether Benedick knows.

Suppose $K_f(K_b(L))$ and $\sim K_{f'}(K_b(L))$.

Florist 1 knows that Benedick knows that Beatrice loves him.

Florist 2 does not know.

So florist 1 will expect a wedding and put down a deposit with a supplier in Napoli while the second florist will not.

They have the same utilities and the same knowledge *about the world* but their knowledge about Benedick's knowledge is different.

And if we know that Beatrice loves Benedick and Benedick knows this, then we can infer the knowledge of the two florists *about Benedick's knowledge* from the fact that one put down a deposit and the other did not.

This is second order reasoning. The difference in the actions of the two florists is not explained by a difference in their knowledge of the world. It is explained by a difference in their knowledge of what Benedick knows.

The first florist knows that proposing to Beatrice is Benedick's dominant strategy. Given that Beatrice loves him she will accept him and there will be a wedding. So preparing to furnish the flowers is his own dominant strategy. But this is not the case for the other florist.

Here is the payoff matrix for either florist. "proposes" refers to whether Benedick proposes.

	proposes	noproposes
order flowers	100	-20
no order	0	0

Thus for florist 1, ordering is the dominant strategy. But not so for the other florist.

But both florists are rational given their state of knowledge.

3 Formalism

We create a language to talk about various knowledge properties in the following way.

- An atomic predicate P is a formula
- If A, B are formulas then so are $\neg A$ and $A \wedge B$
- If A is a formula and i is an agent then $K_i(A)$ is a formula
- We may also include formulas $C(A)$ if we wish to denote common knowledge

3.1 Intuition

Intuitively $K_i(A)$ means that the agent i knows the fact expressed by the formula A . $K_j K_i(A)$ means that j knows that i knows A .

If i, j are the only agents, then $C(A)$ means that i knows A , j knows that i knows A , i knows that j knows that i knows A and so on forever.

For example suppose Ravi and Usha are playing cards. There is a mirror behind Usha so Ravi can see her cards. But Usha does not know this.

And there is a mirror behind Ravi, and Usha can see Ravi's cards. But Ravi only knows about Usha's mirror. And Usha only knows about Ravi's mirror.

Suppose Ravi has the Queen of spades. Let Q represent this fact. Then we have $K_r(Q)$, of course, $K_u(Q)$, and naturally $K_u K_r(Q)$.

But we do not have $K_r K_u(Q)$. Surely this situation will affect the play.

A card game creates an epistemic situation since both players know what cards they have themselves and the default is that the other player does not what cards one has. They play in view of this situation. But the mirrors change the epistemic situation and the game is also played differently. See Dretske [4] who loves card games and uses them in his examples.

3.2 Animal Cognition

3.2.1 Inducing false beliefs in the tigers of the Sundarbans

The Sundarbans are an area at the border of India and Bangladesh where lush forests grow and tiger attacks on humans have been common [18].

Fishermen and bushmen then created masks made to look like faces to wear on the back of their heads because tigers always attack from behind. The payoff matrix for the tiger is below.

	face	noface
attack	-20	1000
not attack	0	0

If the tiger sees a face then his dominant strategy is not to attack since there might be resistance. Thus not attacking is dominant. If the tiger does not see a face then attacking is the dominant strategy.

Thus the fishermen changed the dominant strategy of the tiger by changing its beliefs.⁵ In 1987 no one wearing a mask was killed by a tiger, but 29 people without masks were killed.

Unfortunately the tigers eventually realized it was a hoax, and the attacks resumed.⁶

3.3 The tiger in the bathroom

Suppose I know T , that there is a tiger in your bathroom. I also know that you need to go.

If $\sim K_y(T)$ then you will proceed to the bathroom.

If $K_y(T)$ then you will go the neighbor's apartment and ask if you can use his bathroom. Or perhaps you will call your mother for advice.

⁵I am using the word belief in a weak sense in which we can use it for non-linguistic creatures.

⁶Something which puzzles me is how they passed the knowledge "it is a hoax" from one tiger to another. Tigers are solitary beasts and do not have cellphones.

So I can infer what you know from what you do.

There is one proposition, “a tiger is in the bathroom” which may be true or false and two possible actions for you.

	tiger	no tiger
use own bathroom	-20000	10
neighbor bathroom	-5	-5

The -5 in the second row has to do with the fact that going to the neighbor’s bathroom has a social cost.

If you know about the tiger then the bottom row dominates the top row. If you do not know about the tiger and “no tiger” is your default assumption, then the top row dominates.

What about me? If I see you heading to your own bathroom then I conclude you do not know about the tiger. Under most circumstances my own dominant strategy is to tell you about the tiger. But perhaps I want you to be eaten by the tiger. Then I will not tell you. So my own payoffs are also involved in what I do.

4 A formal framework

Note: we assume that our games are *generic* so that if a player knows what the world is like and what the other players are playing, her best action is uniquely given.

We have n players and some propositions P about the world whose truth value they may or may not know. T is all truth assignments on P .

We define an *epistemic game* with n players to be a map F from T (truth assignments) and S (strategy profiles) to P (payoff profiles)

$$F : T \times S \longrightarrow P$$

So $(F(t, s))_i$ is the payoff to player i when the truth values are according to t and the strategy profile is s .

Note that the payoff does not depend only on the strategy profile but also on the state of the world.

We let $s_i^- = s''$ to mean the strategy profile of all players other than i . We will drop the subscript i when clear from the context.

Let s, s' be strategies for i . we let $s <_t^i s'$ to mean $(\forall s'')(F(t, (s, s'')_i) < F(t, (s', s'')_i))$ (we will usually assume that payoffs for i are never the same so that we need not worry about $<$ and \leq .) In other words we are saying that if t is true then s' is strictly better for i than s no matter what the other players do.

We write $s <_\varphi^i s'$ to mean that for all $t \models \varphi$, $s <_t^i s'$. Again, if φ is true then s' is better for s regardless of which t satisfying φ holds.

We leave out the superscript i on occasion when the context makes it clear.

Theorem 4.1. If $s <_\varphi s'$ and $\psi \models \varphi$ then $s <_\psi s'$

If $s <_\varphi s'$ and $s <_\psi s'$ then $s <_{\varphi \vee \psi} s'$

Corollary 4.2. The set $\{\varphi | s <_\varphi s'\}$ is an *ideal* in the boolean algebra of propositions.

Note that if a rational player knows φ and $s <_\varphi s'$ then the agent will not play s . Moreover if j knows that i knows φ and $s <_\varphi^i s'$ then j knows that the agent i will not play s , and j only needs to respond to strategies other than s . Indeed what j knows about what other players know allows j to reduce the strategy profiles that he needs to respond to.

4.1 Rationalizability

The notion of rationalizability and dominated strategy have been much discussed in the literature [15]. When there is common knowledge of rationality then player i knows that player j will not play a dominated strategy. Given this some of player i 's own strategies can become dominated and i will eliminate them in turn. When this process of elimination of dominated strategies ends, the strategies which remain are the rationalizable ones.

Let us give an example. Suppose player 1 has three strategies a, b, e. Player 2 also has three strategies c, d, f.

a is the best reply to c,f. c is the best reply to b, e. b is the best reply to d. And d is the best reply to a.

So e and f are not best replies to anything and are not rationalizable.

The other four strategies a,b,c,d are rationalizable but there is no (pure) Nash equilibrium.

For instance (a,c) is not a Nash equilibrium because while a is the best reply to c, d is a better reply than c to a.

However, in this situation we have only one payoff matrix known to everyone. But if the payoff matrix depends on the world and different players have different knowledge about the world then the issue becomes complex. But we are confident that the goal described below can be achieved.

Goal To define the notion of rationalizability relative to a given Kripke structure and an epistemic game.

Conjecture: Every strategy rationalizable relative to a Kripke structure is rationalizable in the usual sense. The reverse of course is not true.

Here is a rough argument. Every piece of knowledge you acquire in terms of the world or in terms of what another agent does reduces the possibilities of strategy profiles you face. That means that the relation of dominance becomes larger. So some strategies which might have been rational are so no longer. This happens not only when you yourself learn about the world but also when you learn that someone else has learned about the world, or even about the knowledge of a third agent.

For instance, the strategy of not proposing to Beatrice was rationalizable for Benedick since Beatrice might not love him. But once he knows that she does love him, the strategy of not proposing is not rationalizable. There are *fewer* rationalizable strategies when we know more.

Conjecture: As the agents come to know more, the set of rationalizable strategies decreases.

An *infor* α for agent i is a pair (t, S) where t is a truth assignment and S is an $(n-1)$ -tuple of the other players' strategies.

In a generic game an infor α generates a unique strategy $b(\alpha)$ for agent i which yields the highest value given the infor.

A *state of knowledge* for agent i is a set X of infos.

A strategy s is *rational* for agent i relative to a state of knowledge X if $s \in \{b(\alpha) | \alpha \in X\}$

If $X \subseteq Y$ and s is rational relative to X then it is rational relative to Y .

The less you know, the more rational you are!

A strategy profile (s_1, \dots, s_n) is rational for the agents relative to a tuple of knowledge (X_1, \dots, X_n) if each s_i is rational relative to X_i .

However, not all n-tuples (X_1, \dots, X_n) are possible. For instance if i knows that j knows P then j cannot be playing a strategy t which is dominated when P is true. And i herself cannot be playing a strategy which is dominated when j is not playing t .

So there are connections among the X_i which we have yet to fully investigate.

5 Inferring preferences from actions

1. Example 1

- Jack to Bill: *I am sorry to hear about the fire at your warehouse last night!*
- Bill: *Shhh! It is tomorrow night!*

2. Example 2

- Jack to Bill: *Congratulations on your daughter's wedding yesterday!*
- Bill: *Actually, it is tomorrow.*

3. Example 3

- Jack to Bill: *I hope your wife's surprise birthday party went well!*
- Bill: *Shhh! Actually, it is tomorrow.*

Now why is example 1 a joke and examples 2 and 3 are not? And why do examples 1 and 3 contain a "Shhh!" but example 2 does not?

Examples 2 and 3 are about events which are *good* from the point of view of Bill. But example 1 is presumably bad. But if Bill knows about the fire, why is he not preventing it? We conclude that for Bill

$$fire > \neg fire$$

But if he *prefers* the fire then presumably he is insured and setting the fire or causing it to be set is a crime. Hence the *Shhh!*.

In example 3, the party is a surprise for the wife and that is why Jack is asked to be quiet about it.

6 Conclusion

This paper is work in progress. We have described an approach towards understanding the behavior of groups of agents, by inferring beliefs by watching actions, and by affecting actions by affecting beliefs. Almost all we have said is common sense and the only new thing is to suggest that a formal framework is possible. We have in fact gone a long way towards defining such a framework but some more is to come.

Readers of this paper might now enjoy watching episodes of Adrian Monk, Columbo, or Sherlock Holmes with a new eye and see how the framework applies.

If I may be so daring, even the actions of Obama or Clinton and the way they differ from their statements become explainable.

Acknowledgement Dov Samet very kindly showed me some related work of his [5] which does talk about dominated strategies. This work was done independently of ours and has some elegant ideas. But he and his co-author do not make use of Kripke structures, or, for that matter, detectives and tigers! Thanks to David Makinson for comments. This research was supported by grants from the CUNY Faculty research assistance program.

7 Appendix: Kripke structures

Kripke structures are used to interpret the language above.

Kripke structure M for knowledge for n knowers consists of a space W of states and for each knower i a relation $R_i \subseteq W \times W$.⁷

There is a map π from $W \times A \rightarrow \{0, 1\}$ which decides the truth value of atomic formulas at each state.

We now define the truth values of formulas as follows:

⁷The R_i are often assumed to be equivalence relations and we shall follow this tradition.

1. $M, w \models P$ iff $\pi(w, P) = 1$
2. $M, w \models \neg A$ iff $M, w \not\models A$
3. $M, w \models A \wedge B$ iff $M, w \models A$ and $M, w \models B$
4. $M, w \models K_i(A)$ iff $(\forall t)(wR_it \rightarrow M, t \models A)$

$K_i(A)$ holds at w , (i knows A at w) iff A holds at all states t which are R_i accessible from w .

7.1 Some Consequences

If R_i is reflexive then we will get $K_i(A) \rightarrow A$ (veridicality) as a consequence.

Moreover, regardless of the properties of R_i , we have,

1. If A is logically valid, then A is known
2. If A and $A \rightarrow B$ are known, then so is B

This is the well known problem of *Logical Omniscience* since we are attributing knowledge properties to agents which they do not actually have.

Still, in *small* settings, such assumptions are reasonable.

7.2 Axiom system

1. All tautologies of the propositional calculus
2. $K_i(A \rightarrow B) \rightarrow (K_i(A) \rightarrow K_i(B))$
3. $K_i(A) \rightarrow A$
4. $K_i(A) \rightarrow K_iK_i(A)$
5. $\neg K_i(A) \rightarrow K_i(\neg K_i(A))$

There are also two rules of inference. Modus Ponens, to infer B from A and $A \rightarrow B$. And the other is generalization, to infer $K_i(A)$ from A .

The second rule *does not* say that if A is true than i knows it. Only that it A is a logical truth then i knows it.

These rules are complete. All valid formulas are provable using the axioms and rules. For a bit more detail, see [9].

7.3 Revising Kripke structures when an announcement is made

Suppose we are given a Kripke structure \mathcal{M} . Then some formula φ is announced publicly.

The new Kripke structure is then obtained by deleting all states in \mathcal{M} where φ did not hold. See for instance [13].

References

- [1] Aumann, Robert J. (1976). "Agreeing to Disagree". *The Annals of Statistics* 4 (6): 1236-1239
- [2] Bernheim, D. (1984) Rationalizable Strategic Behavior. *Econometrica* 52: 1007-1028.
- [3] Brandenburger, Adam , *The Language of Game Theory*, World Scientific, 2014.
- [4] Dretske, Fred, *Knowledge and the Flow of Information*, MIT press, (1981).
- [5] Hillas, J. and Dov Samet, "Weak dominance, a mystery cracked", unpublished 2015.
- [6] Kahneman, Daniel, and Amos Tversky. "Prospect theory: An analysis of decision under risk." *Econometrica: Journal of the Econometric Society* (1979): 263-291.
- [7] Lurz, Robert W. *Mindreading animals: the debate over what animals know about other minds*. MIT Press, 2011.

- [8] Lurz, Robert. "If chimpanzees are mindreaders, could behavioral science tell? Toward a solution of the logical problem." *Philosophical Psychology* 22.3 (2009): 305-328.
- [9] Parikh, Rohit. "Recent issues in reasoning about knowledge." *Proceedings of the 3rd conference on Theoretical aspects of reasoning about knowledge*. Morgan Kaufmann Publishers Inc., 1990, pp. 3-10.
- [10] Pearce, D. (1984) Rationalizable Strategic Behavior and the Problem of Perfection. *Econometrica* 52: 1029-1050.
- [11] Parikh, Rohit, and Ramaswamy Ramanujam. "A knowledge based semantics of messages." *Journal of Logic, Language and Information* 12.4 (2003): 453-467.
- [12] Parikh, Rohit, Çağıl Taşdemir, and Andreas Witzel. "The power of knowledge in games." *International Game Theory Review* 15.04 (2013): 1340030.
- [13] Plaza, Jan. "Logics of public communications." *Synthese* 158.2 (2007): 165-179.
- [14] Premack, David, and Guy Woodruff. "Does the chimpanzee have a theory of mind?." *Behavioral and brain sciences* 1.04 (1978): 515-526.
- [15] Rubinstein, Ariel, *Lecture notes in microeconomic theory: the economic agent*. Princeton University Press, 2012.
- [16] Savage, Leonard J. *The foundations of statistics*. Courier Corporation, 1972.
- [17] Shakespeare, William. *Much ado about nothing*. Vol. 12. Classic Books Company, 2001.
- [18] Simons, Marlise, "Face Masks Fool the Bengal Tigers", *The New York Times*, September 5, 1989.
- [19] Tversky, Amos, and Daniel Kahneman. "The framing of decisions and the psychology of choice." *Science* 211.4481 (1981): 453-458.

- [20] Wimmer, Heinz, and Josef Perner. "Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception." *Cognition* 13.1 (1983): 103-128.