

# Perspectives on Bounded Rationality

by: Robert Aumann

The Hebrew University of Jerusalem

Research supported by the National Science Foundation under grant IRI-8814953. Previous versions of this work were presented as the Nancy L. Schwartz Memorial Lecture at Northwestern University in 1986, and at a workshop on Bounded Rationality at the Institute for Mathematical Studies in the Social Sciences (Economics), Stanford University, July, 1989.

## Contents

0. Introduction
1. Evolution
2. Perturbations and Rationality
3. Automata, Computers, Turing Machines
4. Relaxation of Rationality Postulates
5. An Open Problem

## 1. Introduction

Economists have for long expressed dissatisfaction with the complex models of strict rationality that are so pervasive in economic theory. There are several objections to such models. First, casual empiricism or even just simple introspection lead to the conclusion that even in quite simple decision problems, most economic agents are not in fact maximizers, in the sense that they do not scan the choice set and consciously pick a maximal element from it. Second, such maximizations are often quite difficult, and even if they wanted to, most people (including economists and even computer scientists) would be unable to carry them out in practice. Third, polls and laboratory experiments indicate that people often fail to conform to some of the basic assumptions of rational decision theory. Fourth, laboratory experiments indicate that the conclusions of rational analysis (as distinguished from the assumptions) sometimes fail to conform to "reality." And finally, the conclusions of rational analysis sometimes seem unreasonable even on the basis of simple introspection.

From my point of view, the last two of the above objections are more com-

elling than the first three. In science, it is more important that the conclusions be right than that the assumptions sound reasonable. The assumption of a gravitational force seems totally unreasonable on the face of it, yet leads to correct conclusions. “By their fruits ye shall know them” (Matthew).

In the sequel, though, we shall not hew strictly to this line; we shall examine various models that between them, address all the above issues.

To my knowledge, this area was first extensively investigated by Herbert Simon. Much of Simon’s work was conceptual rather than formal. For many years after this initial work, it was recognized that the area was of great importance, but the lack of a formal approach impeded its progress. Particular components of Simon’s ideas, such as satisficing, were formalized by several workers, but never led to an extensive theory, and indeed did not appear to have significant implications that went beyond the formulations themselves.

There is no unified theory of bounded rationality, and probably never will be. Here we examine several different but related approaches to the problem, which have evolved over the last ten or fifteen years. We will not survey the area, but discuss some of the underlying ideas. For clarity, we may sometimes stake out a position in a fashion that is more one-sided and extreme than we really feel; we have the highest respect and admiration for all the scientists whose work we cite, and beg them not to take offense.

From the point of view of the volume of research, the field has “taken off” in the half dozen years. An important factor in making this possible was the development of computer science, complexity theory, and so on, areas of inquiry that created an intellectual climate conducive to the development of the theory of bounded rationality. A significant catalyst was the experimental work of Robert Axelrod in the late seventies and early eighties, in which experts were asked to prepare computer programs for playing the repeated prisoner’s dilemma. The idea of a computer program for playing repeated games presaged some of the central ideas of the later work; and the winner of Axelrod’s tournament —tit-for-tat— was, because of its simplicity, nicely illustrative of the bounded rationality idea. Also, repeated games became the context of much of the subsequent work.

The remainder of these notes is divided into five parts. First we discuss the evolutionary approach to optimization—and specifically to game theory—and some of its implications for the idea of bounded rationality, such as the development of truly dynamic theories of games, and the idea of “rule rationality” (as opposed to “act rationality”). Next comes the area of “trembles”, including equilibrium refinements, “crazy” perturbations, failure of common knowledge of rationality, the limiting average payoff in infinitely repeated games as an expression of bounded rationality, epsilon equilibria, and related topics. Section 3 deals with players who are modelled as computers (finite state automata, Turing machines), which has now become perhaps the most active area in the field. In Section 4 we discuss the work on the foundations of decision theory that deals with various paradoxes (such as Allais and Ellsberg) and with results of laboratory experiments by relaxing various of the postulates and so coming up with a weaker

theory. Section 5 is devoted to one or two open problems.

Most of these notes are set in the framework of non-cooperative game theory, because most of the work has been in this framework. Game theory is indeed particularly appropriate for discussing fundamental ideas in this area, because it is relatively free from special institutional features. The basic ideas are probably applicable to economic contexts that are not game-theoretic (if there are any).

## 2. Evolution

### 2.1 Nash Equilibria as population equilibria

One of the simplest, yet most fundamental ideas in bounded rationality — indeed in game theory as a whole—is that no rationality at all is required to arrive at a Nash equilibrium; insects and even flowers can and do arrive at Nash equilibria, perhaps more reliably than human beings. The Nash equilibria of a strategic (normal) form game correspond precisely to population equilibria of populations that interact in accordance with the rules—and payoffs—of the game.

A version of this idea—the evolutionarily stable equilibrium – was first developed by John Maynard Smith in the early seventies and applied by him to many biological contexts (most of them animal conflicts within a species). But the idea applies also to Nash Equilibria; and not only to interaction within a species, but also to interactions between different species. It is worthwhile to give a more precise statement of this correspondence.

Consider, then, two populations—let us first think of them as different species—whose members interact in some way. It might be predator and prey, or cleaner and host fish, or bees and flowers, or whatever. Each interaction between an individual of population A and one of population B results in an increment (or decrement) in the fitness of each; recall that the *fitness* of an individual is defined as the expected number of its offspring (I use “its” on purpose, since strictly speaking, reproduction must be asexual for this to work). This increment is the payoff to each of the individuals for the encounter in question. The payoff is determined by the genetic endowment of each of the interacting individuals (more or less aggressive or watchful or keen-sighted or cooperative, etc.). Thus one may write a bimatrix in which the rows and columns represent the various possible genetic endowments of the two respective species (or rather those different genetic endowments that are relevant to the kind of interaction being examined), and the entries represent the single encounter payoffs that we just described. If one views this bimatrix as a game, then the Nash equilibria of this game correspond precisely to population equilibria; that is, under asexual reproduction, the proportions of the various genetic endowments within each population remain constant from generation to generation if and only if these proportions constitute a Nash equilibrium.

This is subject to the following qualification: In each generation, there must be at least a very small proportion of EACH kind of genetic endowment; that is, each row and column must be represented by at least SOME individuals. This minimal presence, whose biological interpretation is that it represents possible mutations, is to be thought of as infinitesimal; specifically, an encounter between TWO such mutants (in the two populations) is considered impossible.

A similar story can be told for games with more than two players, and for evolutionary processes other than biological ones; e.g. economic evolution, like the development of the QWERTY typewrite keyboard, about which the economic historian Paul David has written. It also applies to learning processes that are perhaps not strictly analogous to asexual reproduction. And though it does not apply to sexual reproduction, still one may hope that roughly speaking, similar ideas may apply.

One may ask, who are the “players” in this “game”? The answer is that the two “players” are the two populations (i.e., the two species). The individuals are definitely NOT the “players”; if anything, each individual corresponds to the pure strategy representing its genetic endowment (note that there is no sense in which an individual can “choose” its own genetic endowment). More accurately, though, the pure strategies represent kinds of genetic endowment, and not individuals. Individuals indeed play no explicit role in the mathematical model; they are swallowed up in the proportions of the various pure strategies.

Some biologists object to this interpretation, because they see it as implying group or species selection rather than individual selection. The player is not the species, they argue; the individual “acts for its own good,” not the good of the group, or of the population, or of the species. Some even argue that it is the gene (or rather the allele) that “acts for its own good,” not the individual. The point, though, is that NOTHING at all in this model really “acts for its own good”; nobody “chooses” anything. It is the process as a whole that selects the traits. The most we can do is ask what is it that corresponds to the player in the mathematical model, and this is undoubtedly the population.

A question that at first seems puzzling is, what happens in the case of interactions within a species, like animal conflicts for females, etc. Who are the players in the game? If the players are the populations, then this must be a one-person game, since there is only one population. But that doesn’t look right either, and it certainly doesn’t correspond to the biological models of animal conflicts.

The answer is that it is a two-person symmetric game, in which both players correspond to the same population. In this case we look not for just any Nash equilibria, but for symmetric ones only.

## 2.2 Evolutionary Dynamics

The question of developing a “truly” dynamic theory of games has long plagued game theorists and economic theorists. (If I am not mistaken, it is one of

the conceptual problems listed by Kuhn and Tucker in 1953 in the introduction to Volume II of “Contributions to the Theory of Games” —perhaps the last one in that remarkably prophetic list to be successfully solved.) The difficulty is that ordinary rational players have foresight, so they can contemplate all of time from the beginning of play. Thus the situation can be seen as a one-shot game each play of which is actually a long sequence of “stage games,” and then one has lost the dynamic character of the situation.

The evolutionary approach outlined above “solves” this conceptual difficulty by eliminating the foresight. Since the process is mechanical, there is indeed no foresight; no strategies for playing the repeated game are available to the “players”.

And indeed, a fascinating dynamic theory does emerge. A contribution to this theory has been made by Young, about which we will hear later today. A book on the subject has been written by Siegmund, and there is an excellent chapter on evolutionary dynamics in the book by van Damme on refinements of Nash equilibrium. Notable recent contributions have been made by Kandori (Princeton), Mailath (Penn) and Rob (Penn). Many others have also contributed to the subject.

It turns out that Nash equilibria are often unstable, and one gets various kinds of cycling effects. Sometimes the cycles are “around” the equilibrium, like in “matching pennies,” but at other times one gets more complicated behavior. For example, the game

0, 0	4, 5	5, 4
5, 4	0, 0	4, 5
4, 5	5, 4	0, 0

has  $((1/3, 1/3, 1/3), (1/3, 1/3, 1/3))$  as its only Nash equilibrium; the evolutionary dynamics does not cycle “around” this point, but rather confines itself (more or less) to the strategy pairs in which the payoff is 4 or 5. This suggests a possible connection with correlated equilibria; this possibility has recently been investigated by Dean Foster (Chicago) and collaborators.

Thus evolutionary dynamics emerges as a form of rationality that is bounded in that foresight is eliminated.

### 2.3 “Rule Rationality” vs. “Act Rationality”

In a famous experiment conducted by Werner Guth (Frankfurt) and collaborators —and later repeated, with important variations, by Ken Binmore (Michigan)—two players were asked to divide a considerable sum of money (ranging as high as DM 100). The procedure was that P1 made an offer, which could be either accepted or rejected by P2; if it was rejected, nobody got anything. The players did not know each other and never saw each other; communication was a one-time affair via computer.

“Rational” play would predict a 99-1 split, or 95-5 at the outside. Yet in by far the most trials, the offered split was between 50-50 and 65-35. This is surprising

enough in itself. But even more surprising is that in most (all?) cases in which P2 was offered less than 30%, he actually REFUSED. Thus he PREFERRED to walk away from as much as DM 25 or 30. How can this be reconciled with ordinary notions of utility maximization, not to speak of game theory?

It is tempting to answer that a player who is offered five or ten percent is “insulted.” Therefore his utilities change; he gets positive probability from “punishing” the other player.

That’s alright as far as it goes, but it doesn’t go very far; it doesn’t explain very much. The “insult” is treated as exogenous. But obviously the “insult” arose from the situation. Shouldn’t we treat the “insult” itself endogenously, somehow explain IT game-theoretically?

I think that a better way of explaining the phenomenon is as follows: Ordinary people do not behave in a consciously rational way in their day to day activities. Rather, they evolve “rules of thumb” that work in general, by an evolutionary process like that discussed at 2.1 above, or a learning process with similar properties. Such “rules of thumb” are like genes (or rather alleles). If they work well, they are fruitful and multiply; if they work poorly, they become rare and eventually extinct.

One such rule of thumb is “Don’t be a sucker; don’t let people walk all over you.” In general, the rule works well, so it becomes widely adopted. As it happens, the rule doesn’t apply to Guth’s game, because in that particular situation, a player who refuses DM 30 does not build up his reputation by the refusal (because of the built-in anonymity). But the rule has not been consciously chosen, and will not be consciously abandoned.

So we see that the evolutionary paradigm yields a third form of bounded rationality: Complicated strategies that require a lot of computation are replaced by simple rules of thumb that work well “on the whole.”

### 3. Perturbations of Rationality

#### 3.1 Equilibrium Refinements

Equilibrium refinements (Selten, Myerson, Kreps, Wilson, Kalai, Samet, Kohlberg, Mertens, Weibull, Van Damme, Reny, Cho and many others) don’t really sound like bounded rationality. They sound more like super-rationality, since they go beyond the basic utility maximization that is inherent in Nash equilibrium. In addition to Nash equilibrium, which demands rationality on the equilibrium path, they demand rationality also off the equilibrium path. Yet all are based in one way or another on “trembles”—small departures from rationality.

The paradox is resolved by noting that in a game situation, one man’s ir-

rationality requires another one's superrationality. YOU must be superrational in order to deal with MY irrationalities. Since this applies to all players, taking account of possible irrationalities leads to a kind of superrationality for all. To be superrational, one must leave the equilibrium path. Thus a more refined concept of rationality cannot feed on itself only; it can only be defined in the context of irrationality.

### 3.2 Crazy Perturbations

An idea related to the trembling hand is the theory of irrational or "crazy" types, as propounded first by the "gang of four" (Kreps, Milgrom, Roberts, and Wilson), and then taken up by Fudenberg, Maskin, Aumann and Sorin, Levine, Gilboa and Samet, and no doubt others. In this work there is some kind of repeated or other dynamic game set-up; it is assumed that with high probability the players are "rational" in the sense of being utility maximizers, but that with a small probability, one or both play some one strategy, or one of a specified set of strategies, that are "crazy"—have no a priori relationship to rationality. An interesting aspect of this work, which differentiates it from the "refinement" literature, and makes it particularly relevant to the theory of bounded rationality, is that it is usually the crazy type, or a crazy type, that wins out – takes over the game, so to speak. Thus in the original work of the gang of four on the prisoner's dilemma, there is only one crazy type, who always plays tit-for-tat, no matter what the other player does; and it turns out that the rational type must imitate the crazy type, he must also play tit-for-tat, or something quite close to it. Also, the "crazy" types, while irrational in the sense that they do not maximize utility, are usually by no means random or arbitrary (as they are in refinement theory). For example, we have already noted that tit-for-tat is computationally a very simple object, far from random. In the work of Aumann and Sorin, the crazy types are identified with bounded recall strategies; and in the work of Fudenberg and Levine, the crazy types form a denumerable set, suggesting that they might be generated in some systematic manner, e.g. by Turing machines. There must be method to the madness; this is associated with computational simplicity, which is another one of the underlying ideas of bounded rationality.

### 3.3 Epsilon-equilibria

Rather than playing irrationally with a small probability, as in 2.1 and 2.2 above, one may deviate slightly from rationality by playing so as almost, but not quite, to maximize utility; i.e., by playing to obtain a payoff that is within epsilon of the optimum payoff. This idea was introduced by Radner in the context of repeated games, in particular of the repeated prisoner's dilemma; he showed that in a long but finitely repeated prisoner's dilemma, there are epsilon-equilibria with small epsilon in which the players "cooperate" until close to the end (though, as

is well-known, all exact equilibria lead to a constant stream of “defections”).

### 3.4 Infinitely Repeated Games with Limit-of-the-Average Payoff

There is an interesting connection between epsilon-equilibria in finitely repeated games and infinitely repeated games with limit of the average payoff (“undiscounted”). The limit of the average payoff has been criticized as not representing any economic reality; many workers prefer to use either the finitely repeated game or limits of equilibrium payoffs in discounted games with small discounts. Radner, Maskin, Myerson, Fudenberg, Mertens, Neyman, Forges and perhaps others have demonstrated that the results of these two kinds of analysis can indeed be quite different.

Actually, though, the infinitely repeated undiscounted game is in some ways a simpler and more natural object than the discounted or finite games. In calculating equilibria of a finite or discounted game, one must usually specify the number  $n$  of repetitions or the discount rate  $d$ ; the equilibria themselves depend crucially on these parameters. But one may want to think of such a game simply as “long”, without specifying how long. Equilibria in the undiscounted game may be thought of as “rules of thumb,” which tell a player how to play in a “long repetition, independently of HOW long the repetition is. Whereas limits of finite or discounted equilibrium payoffs tell the players approximately how much payoff to expect in a long repetition, analysis of the undiscounted game tells him approximately how to play.

Thus the undiscounted game is a framework for formulating the idea of a duration-independent strategy in a repeated game. Indeed, it may be shown that an equilibrium in the undiscounted game is an approximate equilibrium simultaneously in all the  $n$ -stage truncations, the approximation getting better and better all the time. Formally, a strategy profile (“tuple”) is an equilibrium in the undiscounted game if and only if for some sequence of  $\epsilon(n)$  tending to zero, each of its  $n$ -stage truncations is an  $\epsilon(n)$ -equilibrium (in the sense of Radner described above) in the  $n$ -stage truncation of the game.

### 3.5 Failure of Common Knowledge of Rationality

In their paper on the repeated prisoner’s dilemma, the Gang of Four pointed out that the effect that they were demonstrating holds not only when one of the players believes that with some small probability, the other is a tit-for-tat automaton, but also if one of them only believes (with small probability) that the other believes this about him (with small probability). More generally, it can be shown that many of the perturbation effects we have been discussing do not require an actual departure from rationality on the part of the players, but only a lack of common knowledge of rationality.

#### 4. Automata, Computers and Turing Machines

We come now to what is probably the “mainstream” of the newer work in bounded rationality, namely the theoretical work that has been done in the last four or five years on automata and Turing machines playing repeated games. The work was pioneered by A. Neyman and A. Rubinstein, working independently and in very different directions. Subsequently, the theme was taken up by E. Ben-Porath, E. Kalai, W. Stanford, E. Zemel, D. Abreu, B. Peleg, E. Lehrer, A. Mor, C. Papadimitriou, R. Stearns, E. Einy, and many others, each of whom made significant new contributions to the subject in various different directions. Different branches of this work has been started by A. Lewis and K. Binmore, who have also had their following.

It is impossible to do justice to all this work in a reasonable amount of space, and we content ourselves with brief descriptions of some of the major strands. In one strand, pioneered by Neyman, the players of a repeated game are limited to using mixtures of pure strategies each of which can be programmed on a finite automaton with an exogenously fixed number of states. This is reminiscent of the work of Axelrod, who required the entrants in his experiment to write the strategies in a fortran program not exceeding a stated limit in length. In another strand, pioneered by Rubinstein, the size of the automaton is endogenous; computer capacity, so to speak, is considered costly, and any capacity that is not actually used in equilibrium play is discarded. The two approaches lead to very different results. The reason is that Rubinstein’s approach precludes the use of “punishment” or “trigger” strategies, which swing into action only when a player departs from equilibrium, and whose sole function is precisely to prevent such departures. In the evolutionary interpretation of repeated games, Rubinstein’s approach may be more appropriate when the stages of the repeated game represent successive generations, whereas Neyman’s may be more appropriate when each generation plays the entire repeated game (which would lead to the evolution of traits having to do with reputation, like “Don’t be a sucker”).

The complexity of computing an optimal strategy in a repeated game, or even just a best response to a given strategy, has been the subject of works by several authors, including Ben Porath and Papadimitriou. Related work has been done by Alain Lewis, though in the framework of recursive function theory (which is related to infinite Turing machines) rather than complexity theory (which has to do with finite computing devices). Roughly speaking, the results are qualitatively similar: finding maxima is hard. Needless to say, in the evolutionary approach to games, nobody has to find the maxima; they are picked out by evolution. Thus the results of complexity theory again underscore the importance of the evolutionary approach.

Binmore and his followers have modelled games as pairs (or  $n$ -tuples) of Turing machines in which each machine carries in it some kind of idea of what the other “player” (machine) might look like.

Other important strands include work by computer scientists who have made

the connection between distributed computing and games (“computers as players”, rather than “players as computers”). We will not venture to survey this extensive work, about which the audience here knows a lot more than us.

## 5. Relaxation of Rationality Postulates

A not uncommon activity of decision, game, and economic theorists since the fifties has been to call attention to the strength of various postulates of rationality, and to investigate the consequences of relaxing them. Many workers in the field — including the writer of these lines—have at one time or another done this kind of thing. People have constructed theories of choice without transitivity, without completeness, violating the sure-thing principle, and so on. Even general equilibrium theorists have engaged in this activity, which may be considered a form of limited rationality (on the part of the agents in the model). This kind of work is most interesting when it leads to outcomes that are qualitatively different—not just weaker—from those obtained with the stronger assumptions; but I don’t recall many such cases. It can also be very interesting and worthwhile when one gets roughly similar results with significantly weaker assumptions.

## 6. An Open Problem

We content ourselves with one open problem, which is perhaps the most challenging conceptual problem in the area today: to develop a meaningful formal definition of rationality in a situation in which calculation and analysis themselves are costly and/or limited. In the models we have discussed up to now, the problem has always been well defined, in the sense that an absolute maximum is chosen from among the set of feasible alternatives, no matter how complex a process that maximization may be. The alternatives themselves involve bounded rationality, but the process of choosing them does not.

Here, too, an evolutionary approach may eventually turn out to be the key to a general solution.