

# Bounded Rationality :: Bounded Models

**Jocelyn Smith**

University of British Columbia  
201-2366 Main Mall  
Vancouver BC  
jdsmith@cs.ubc.ca

## Abstract

In economics and game theory agents are assumed to follow a model of perfect rationality. This model of rationality assumes that the rational agent knows all and will take the action that maximizes her utility. We can find evidence in the psychology and economics literature as well as in our own lives that shows human beings do not satisfy this definition of rationality. Thus there are many who look to study some notion of bounded rationality. Unfortunately, models of bounded rationality suffer from the exact phenomena that they attempt to explain. Specifically, models of bounded rationality are bounded. Understanding the limits of various rationality models will make clearer their contribution and place in the overall picture of rationality.

## Introduction

There are three primary motivations for defining and studying bounded rationality. The first is that under the assumption of perfect rationality, "rational decision can be detrimental to the satisfaction of self-interest" -Max Black. The finitely repeated prisoner's dilemma, which will be presented in detail later, is the archetypal example of this less than satisfactory situation. The second reason for studying bounded rationality is the computational infeasibility of perfect rationality in all but the simplest of situations. The typical example here is the game of chess though any NP-hard problem will suffice. The third reason for being interested in bounded rationality is that humans rarely exemplify the perfect rationality model. Examples of this can be found in psychology and experimental economics experiment results. These deviations include choosing differently among the same two options if they are presented in different by equivalent contexts (framing effects), simplifying even basic problems and thus picking suboptimal actions, and a lack of indifference to irrelevant alternatives.

Any paper dealing with bounded rationality will cite some form of one or all of these as motivation for the work. However, it is often the case that the bounded rationality model(s) that follow these motivations are not linked back to them. Thus the reader is left to determine how and when the given model actually addresses these problems with perfect rationality. In most cases, the model addresses only one of these

issues and sometimes only in very specific instances. It is often not obvious which issue a particular model addresses and what that model's contribution is in the larger context of understanding bounded rationality. In this paper, we present several models from the literature and determine which of these issues they address. We then present ideas of how each contributes to our overall understanding of rationality.

## History

Concepts of rationality that deviate from that of perfect rationality are not new. In fact, philosophers have long struggled with and debated the definition of rationality. In his discussion of rationality, Black presents several definitions present by past and present philosophers (Black 1990). For some, rationality is that which differentiates humans from other "lesser" animals. "A faculty in man, that faculty whereby man is supposed to be distinguished from beasts, and wherein it is evident that he much surpasses them." - John Locke (1894). The classical conception of "reason" was as the separate faculty that bridles forces of passion. (note: modern philosophy also regards active intelligence as being motivated by passion) Bertrand Russell defined the rational choice as the choice of the right means for the ends. On the other hand, Micheal Oakeshott regarded behaviour as rational if it was faithful to the individual's knowledge of how to behave well. Thus while some philosophers would define rationality using a definition similar to that of perfect rationality (ex. Bertrand Russell) there are many conflicting views even at this high level. The length of this debate and the variety of definitions highlight the fact that understanding rationality is not a new or simple task.

## Abstract Models

Herbert Simon was a pioneer and verbal proponent of bounded rationality in Economics. His description of bounded rationality is perhaps the closest to human rationality. His theory is based on the idea that there are two approaches to solving intractable problems. The first is optimizing. A boundedly rational agent that optimizes first finds an appropriate approximation of the given problem. The agent then determines the optimum solution to the approximation and uses this as the solution to the initial problem (Simon 1982). We do this anytime we model a really world

problem. Take for example the use of the perfect rationality model in economics. This model is used to find optimal equilibrium and the results are then used to predict real world activities even though the real world situation may not apply perfect rationality. The second strategy is to satisfice. Here the agent will search possible solutions until a solution that is satisfactory is found. For example when searching for a course project topic a student will search until he finds one that is good enough. Note searching all possible topics would prevent the student from ever actually doing a project.

Simon provides strong and convincing arguments that the pursuit of bounded rationality modes is a worthwhile endeavour. His main contribution is the presentation of a framework from which to view bounded rationality. It is important to have such a framework to guide the development of models. It is also important to be aware of the limitations of the framework. One issue to be aware of is that these two strategies may be intermixed. With this in mind the question becomes how can we apply these notions to describe or predict/prescribe rational behaviour.

### Concrete Models

With the knowledge from philosophy that we need to consider what we mean by bounded rationality and Simon's framework as a possible guide we will now look at some more concrete models of bounded rationality.

#### Rational Rules

One of the results of psychology experiments that shows people do not always behave perfectly rationally is that we tend to simplify a problem (even a relatively simple one) and this simplification can lead to choices that do not maximize utility (Rubinstein 1998). One form simplification can take is the cancelling out of similar parts. The motivating example found in Rubinstein is a choice between two lotteries. A lottery where  $x$  is awarded with probability  $p$  and 0 is rewarded with probability  $1 - p$  is represented by  $(x, p)$ .

Given the choice between

$$L_3 = (4000, 0.2) \text{ and } L_4 = (3000, 0.25)$$

the popular choice is  $L_3$ .

However, given the choice between

$$L_1 = (4000, 0.8) \text{ and } L_2 = (3000, 1.0)$$

the most common choice is  $L_2$ .

Notice, however, that  $L_3$  and  $L_4$  can also be expressed as

$$L_3 = 0.25L_1 + 0.75[0] \text{ and } L_4 = 0.25L_2 + 0.75[0]$$

Thus the preference for  $L_3$  on one hand and  $L_2$  on the other violates the von Neumann-Morgenstern independence axiom. The theory is that when choosing between lotteries of these types people try to simplify by cancelling similar

parts of the vectors. So when confronted  $L_3$  and  $L_4$  people determine that 0.2 and 0.25 are similar but 4000 and 3000 are not so they choose  $L_3$ . However, neither part of  $L_1$  and  $L_2$  is considered similar so another choice function, perhaps risk aversion, operates instead. Rubinstein takes this choice rule, formalizes it and shows that there is a preference relation that is consistent with it.

This rule is obviously motivated by the fact that humans are not perfectly rational and it is believable that humans do use rules such as this in their reasoning. The argument for this one is quite convincing. Though Rubinstein presents an argument that this similarity procedure is consistent with rational behaviour he does not go further and explain why this is interesting. Alone this simplification rule may not be that interesting as it simply gives an explanation for human behaviour in a very particular case, choosing between two choices represented as probability vectors. In fact, in most situations, this rule alone is unlikely to predict human behaviour and even in the example used to motivate it there is the assumption of a risk aversion rule being applied when simplification is not sufficient.

The contribution then is that the formalization and argument that this formalization is consistent with rationality given an appropriate preference relation shows a possible way of describing and analyzing rules such as this one as rational. Since the really benefit of understanding such a rule can only be fully realized by understanding it in the context of a set of such rules, there are several interesting questions left. What would a full set or even a substantial subset of such rules look like? When do individual rules apply and how do rules interact? Is it possible to formalize sets of rules and construct a preference relation that works for all of them? Rubinstein's work suggests a way in which experiment knowledge of human behaviour might be used to construct a more formal model of these rational processes.

#### Bounding information

One way to bound rationality is to put constraints or add costs to the information that is used to make the rational decision. Rubinstein identifies three costs of information: acquisition, memory and communicating. The theory presented assumes that the agent optimizes on what to know. However, this can often be a harder optimization problem than the original and thus seems a bit hard to justify in terms of providing understanding of how humans rationalize. Rubinstein claims there are situations where this makes sense but does not back up this claim beyond saying that these are situations where the decision maker regularly applies a rule to make a decision and only occasionally considers which rule to apply taking into account its complexity. One case where this may apply is in terms of developing AI systems where trade offs could be determined while designing the system. In general, however, expecting that even occasionally a hard problem is solved may not be reasonable. Especially in the very likely case that it is an NP-hard optimization problem. Rather the true value of a model such as this may lie in the argument often applied in economics to justify the perfect rationality model. That is that there are some situations in which people tend to act "as if" they are

optimizing the knowledge they acquire, store and communicate taking into account the complexity of doing so. Perhaps using this model could lead to more satisfying equilibria in problems where the perfect rationality model leaves us with only unsatisfying equilibrium.

### A bounded memory model

One bound on rationality is memory. How much information can be stored? One way to model memory is to use a finite automaton and consider the number of states to represent the amount of memory available. In this model, the states represent the agents strategy. States correspond to the agent's actions and transitions correspond to the opponent's actions. This model leads to some interesting results specifically in terms of the finitely repeated prisoner's dilemma game.

**Prisoner's dilemma** The archetypal example of a situation in which perfect rationality leads to unsatisfactory results is the game known as prisoner's dilemma. (Figure1) The story that goes with the game is that two prisoner's are being interrogated by the police. They both have the choice of not revealing anything (cooperating with the other prisoner) or of signing a deal and agreeing to testify against the other (defecting). If both cooperate they will both go to jail for a short period of time on a lesser charge. If one cooperates and the other defects, the defector will go free and the cooperator will serve a long sentence. If both defect, both will go to jail for an intermediary length of time. In prisoner's dilemma the dominate strategy of both players is to defect. This is because no matter what the "column" player does "row" is better off defecting. If column cooperates then row gets a reward of 4 instead of the 3 he would have got if he cooperated. If column defects then row gets a reward of 1, instead of the 0 he would have got if he had cooperated. Since both players will defect, they will both receive a reward of 1. However, it is obviously preferable for both players to cooperate thus both getting a reward of 3.

	C	D
C	(3,3)	(0,4)
D	(4,0)	(1,1)

Figure 1: The payoffs for the prisoner's dilemma.

This situation where both players defect is also the only Nash-equilibrium if the game is played n-times. Note a Nash-equilibrium is a set of strategies such that no individual player can improve his reward by deviating from the given strategy. The proof is by backwards induction. In the last round the best thing to do is to defect (by the same ar-

gument as for the one shot game). Since defection is the only option in the last round there is no reason to collaborate in the second last round and so on. This result seems even less "rational" than that of the one shot game. Note that two versions of prisoner's dilemma that do lead to cooperative play are the infinitely repeated game and the finitely repeated game where the number of games is unknown.

**Automaton with fewer than n-states** By limiting agents to strategies that can be modelled using finite automaton of particular sizes dependent on the number of rounds it is possible to get more satisfying equilibrium (Papadimitriou & Yannakakis 1994). In the n-prisoner's dilemma, if both agents strategies are limited to less than n-states then they can not count to n and so cannot apply backwards induction. In this case, a basic strategy known as tit-for-tat (Figure2) can be shown to be an equilibrium (Papadimitriou & Yannakakis 1994).

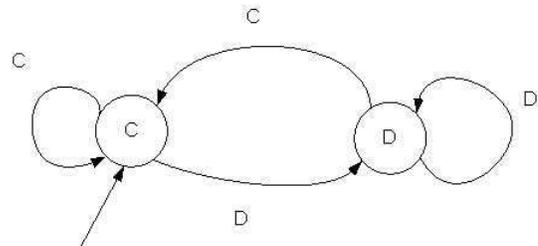


Figure 2: Tit-for-Tat: In repeated prisoner's dilemma, a strategy of always co-operating unless the other defects. If the other defects then defect one round to punish. This strategy can be played using the above two state automaton.

**Automaton with fewer than exponential states** The more interesting case is where the automaton have more than n-states but where at least one has less then exponential in n states.

**Theorem 1** (Papadimitriou & Yannakakis 1994) For every  $\epsilon > 0$ , in the n-round prisoner's dilemma played by automata with size bounds. If at least one of the bounds is  $< 2^{c_\epsilon n}$  then there is a mixed strategy equilibrium with average payoff for each player at least  $3 - \epsilon$ . ( $c_\epsilon = \epsilon/6(1 + \epsilon)$ )

In other words, that it is possible to construct a mixed strategy equilibrium such that both players get payoffs arbitrarily close to the payoff they would receive if they cooperated in every round. The equilibrium strategy involves several looping components. For the first  $d$  rounds the players state a business card or series of moves (selected according to their mixed strategy). Here  $d$  depends on the size bound of the players. Later an agent is required to remember the business card of his opponent and thus must store all  $2^d$  possible cards. After this exchange, the agents play according to some rules which even out any advantage one might have gained over the other in the first  $d$  rounds. Next they loop over two stages. The first is collaboration for a large number of rounds (determined by their size). The second is a sequence that depends on both the presented business cards.

This stage has two purposes: to even out advantageous business cards and to show that agents remember each other's business cards. If ever the other player does not play as expected then the agent goes into a state of perpetual defection. The construction is such that the players must use all states to play the strategy and thus do not have any states remaining to count and apply backward induction. In fact, even if one player does have exponentially many states and can count then the only difference will be that he defects in the last round (since doing so earlier will not increase his payoff).

**Implications** When presenting this work the author's motivate it by saying that it remedies the "unnatural" or even "irrational" predictions of game theory in the n-round prisoner's dilemma. This is true in that by using this model the given strategies are an equilibrium where cooperation prevails. However, this model does not satisfy our desire to model human behaviour for a couple of reasons. Specifically, cooperation that arises in prisoner's dilemma when played by humans is unlikely to arise from such an elaborate equilibrium. Also given a small number of games (say 10) it does not seem satisfactory to say that cooperation arises because we do not have enough memory to count. Especially, considering that knowing the backwards induction argument does not prevent people from collaborating.

However, there are actually several other ways that this work contributes to the bounded rationality literature. One is the "as if" argument used to defend the perfect rationality assumption. That is that people behave "as if", as if they were perfectly rational or in this case as if they had less than exponential memory in the number of rounds. Of course it is unlikely that this is always true anymore than it is always true for the perfect rationality assumption; however, it may be a useful approach when the assumption holds. The other contribution of this work is that it shows how simple (tit-for-tat) and/or complex (sub-exponential equilibria) automaton representation of strategies may be. Finally, perhaps by contrasting the strategies employed by humans with those that invoke equilibria in finite automaton it is possible to learn more about human rationality. For example, that we play more like automaton with less than n-states than like those with more than n-states.

### Limited foresight reasoning

The motivating example for limited foresight reasoning is the game of chess. Each game is a series of moves ending in one of three outcomes. Using a backward induction argument on terminal histories, Zermelo proved at the beginning of the twentieth century that chess has an inherent value (white wins, black wins, or draw). However, this value is unknown since finding it would require enumerating all possible game histories a low estimate of which is  $10^{40}$ . Furthermore, because of this intractability a player can not select the move guaranteed to give her the best possible outcome. Yet without being able to see every possible outcome of the game chess players are able to play "rationally".

Players of chess employ several strategies for determining good moves. One is to select reasonable candidate moves (likely based on past experience and knowledge of the game)

and then to look ahead several moves to see how the game may progress.

A couple simple models of this look ahead behaviour are discussed by Rubinstein. In both models it is assumed that the player knows the equilibrium actions of both himself and his opponent for the next k moves (Rubinstein 1998). The goal for the agent is to maximize the payoffs collected over these k moves. In one model, the agent may only change the current move but not his future moves. In the other, the opponent's actions are assumed to be fixed but the agent can optimize over changes to all of his next k moves. Both models lead to reasonable definitions of equilibrium. However, intuitively they are not very satisfactory. The first because an agent must treat his future moves as determined even though he can influence them. In both the opponents moves are treated as given, though it is likely the opponent's strategy will change depending on how the agent acts. The value of these models is most likely in the fact that they do not provide a very satisfactory description of limited foresight reasoning and thus present the challenge of creating a better model.

### Bounded Optimality

Russell et al. define a model of rationality they believe is appropriate for determining if an AI system exhibits rational behaviour (Russell & Subramanian 1993). To motivate the approach taken four possible definitions of rationality are given. The first three are currently used in AI and the fourth is the proposed definition upon which the paper is based.

1. Perfect rationality: The rational agent always acts in such a way as to maximize his utility.
2. Calculative rationality: The rational agent eventually gives the answer that would have been optimal at the start of its deliberation.
3. Metalevel rationality: The rational agent computes the best sequence of computation plus action given that the computation must select the action.
4. Bounded optimality: Specifies optimal programs. A rational agent is an agent that does as well as possible given its computational resources.

The desired property of intelligent systems was traditionally definition 1 but this is not always computationally feasible. Definition 2 is interesting in so far as it may demonstrate that a system in principle would do the right thing but is not generally useful in practise since it may take exponential time to compute a solution. Finally, definition 3 often presents a problem that is harder than the original. One common approach in AI is to design calculatively rational systems and then use speedups and approximations with the goal of getting close to the optimal. Approximations of meta-level rationality have also proven successful in many applications.

An agent function is a mapping from an agent's perceptions to its actions. An agent's program is what generates its actions. The agent function is defined in (Russell & Subramanian 1993) such that the action may be a null action if the agent is still calculating. Russell et al. construct a

framework for designing bounded optimal agents in specific environments. A bounded optimal agent is defined as an agent that maximizes utility over all functions that are feasible given the time constraints and a specified class of machines. They extend the notion of bounded optimal agents to asymptotically bounded agents where the notion of asymptotic is symmetric to that found in complexity.

Thus this paper formalizes ideas about what is rational when constraints do not permit perfect rationality. The model is not useful in describing human rationality. However, it presents an plausible way for defining boundedly rational AI systems.

## Conclusion

The motivations for developing models of bounded rationality are that the perfect rationality assumption does not always give reasonable results, may be infeasible in practice and does not always reflect human rational behaviour. However, models of bounded rationality are also limited in their ability to address these issues. Black sums up the complexity of the issue, "... in philosophical discussions of rationality, there is a sense in which we do not 'know what we are talking about' and can never do so."(Black 1990)

The complexity of understanding rationality does not mean that it is not worth pursuing the issue. Rather it suggests that any comprehensive theory of rationality must include several models including perfect rationality. Various models make different contributions to our understanding of rationality and are applicable in different situations.

In this paper we have examined several models that address different aspects of bounded rationality. One model that we did not examine but that has contributed significantly to research involving rationality is the model of perfect rationality. One contribution of this model is the economics and game theory that has been built up around it. The areas where this model falls short are exactly those in which a model of bounded rationality may be interesting to examine. Together the models discussed in this paper start to fill in these holes. A high level and fairly abstract model of bounded rationality is discussed by Simon. His work helps to motivate a need for bounded rationality models and suggests two boundedly rational approaches: optimizing an approximation to the problem and satisficing or finding a satisfactory solution to the problem. A much lower level look at human rationality is provided by specific rules such as the simplifying rule formalized by Rubinstein. An approach that considers the cost of information to an agent suggests optimizing over actions considering the complexity of the information required for this action. This meta-level reasoning may be interesting in situations where on a regular basis a rule is applied to make a decision but the decision of which rule to apply need only be made once. Finite automaton models provide another way to study strategies and equilibrium. Interesting comparisons may be made between equilibrium found in practice during human interactions and those possible on finite automaton. Two limited by specific models of optimizing subproblems are found in the k-look-ahead models. Finally, a theory for discussing the optimality of a bounded agent in AI provides an interesting take on

bounded rationality that is motivated by computational complexity. Together these models begin to describe some of the complexity of rationality.

## References

- Black, M. 1990. *Perplexities: Rational Choice, the Prisoner's Dilemma, Metaphor, Poetic Ambiguity, and Other Puzzles*. Cornell University Press.
- Papadimitriou, C. H., and Yannakakis, M. 1994. On complexity as bounded rationality (extended abstract). *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing* 726–733.
- Rubinstein, A. 1998. *Modeling Bounded Rationality*. MIT Press.
- Russell, S., and Subramanian, D. 1993. Provably bounded-optimal agents. *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*.
- Simon, H. 1982. *Models of Bounded Rationality: Behavioral Economics and Business Organization*. MIT Press.