

UC Irvine

UC Irvine Previously Published Works

Title

The Social Contract, the Game of Life and the Shadow of the Future

Permalink

<https://escholarship.org/uc/item/9531t22z>

Author

Skyrms, Brian

Publication Date

2022

DOI

10.1007/s41412-022-00119-6

Peer reviewed

The Social Contract, The Game of Life and The Shadow of the Future

Brian Skyrms

Logic and Philosophy of Science

University of California, Irvine

Irvine, California, USA

Abstract: Ken Binmore's treatment of his Game of Life as a bargaining game and his treatment of morality as an equilibrium selection device for that game, are examined in the context of repeated games with both infinite and finite horizon. With a finite horizon, there are three different viable approaches. They differ in the way they impact his treatment of morality.

Keywords: Social Contract, Game of Life, Bargaining, Shadow of the Future

JEL Codes B5, C7, D6, Y8

1 The Game of Life In social life, there are all sorts of interactions that can be modeled as all sorts of games. This conglomerate is Ken Binmore's

“Game of Life”. Several philosophers have paid special attention to interactions that may be modeled as Nash bargaining games. Binmore goes further, and argues that the Game of Life should be viewed as one grand noncooperative Nash bargaining game. This characterization of the Game of Life relies on the shadow of the future -- that is, on the theory of repeated games. The Folk Theorems of repeated game theory give the payoff structure of a Nash bargaining game, with an infinite number of Pareto-efficient Nash equilibria. How people are to live in social interaction is a question of how to select among this infinite number of equilibria.

Binmore sees ethics as a way of solving this equilibrium selection problem, of coordinating on a Pareto-efficient Nash equilibrium in social life. He goes on to develop his own theory of ethics – the Game of Morals – which is closely related to the veil of ignorance approaches of John Harsanyi and of John Rawls. (There are, of course, differences between Harsanyi and Rawls, which we do not discuss here.) In Binmore’s Game of Morals, individuals move behind a veil of ignorance and choose a solution to the bargaining game, under the assumption that in the Game of Life they have an equal chance of being anyone in society. We are here interested in reductions of the Game of Life to a Nash bargaining game via Folk Theorems, and the perspectives on the Game of Morals that they supply.

These original Folk Theorems bring in infinity in a way that is technically nice, but ultimately hard to justify. Repetitions of a game cannot have a finite bound. Either the game is repeated infinitely often, or there is fixed positive probability that it will be repeated again no matter how many times it has already been played. Strategies in the repeated game then can support behavior that would be unavailable in single plays of the game. A notable example is cooperation in Prisoner's Dilemma. But if the game is repeated a finite number of times, no matter how great, the theorems fail. For finitely repeated Prisoner's Dilemma, the only Nash equilibrium is "Always Defect."

There are, however, relevant approaches with finite repeated games. More positive results are available for three main cases. (1) Agents with bounded reasoning capacity can only consider simple strategies in the repeated game. This restriction of eligible strategies may remove those responsible for cooperation unraveling in long finitely repeated games. (2) Agents may settle for strategies whose payoffs are close to the optimal. This "epsilon best response" allows for corresponding epsilon-equilibria, and thus expands the set of equilibria beyond Nash equilibria. Some of these equilibria can sustain cooperation. (3) Players can be uncertain about other players' types, as in Harsanyi's games of incomplete information. This shifts the focus to players beliefs, and the relevant equilibrium concept to Bayes-Nash equilibrium. There are Bayes-Nash equilibria that support considerable

cooperation in finitely repeated Prisoner's Dilemma. These three approaches have been developed in ways that show how the reduction can be maintained. The interpretation, however, is changed in ways relevant to social theory in general, and the Game of Morals in particular.

The next section revisits the genesis of repeated game theory in the early days of game theory at the RAND corporation and shows that the relevant ideas for cooperation in finitely repeated games were already there. Then we give a brief overview of the main relevant results first, for infinitely repeated games with discounting of the future and second, for finitely repeated games. Finally, we consider the import for the Game of Morals of moving to finitely repeated games.

2 John Nash and The Prisoner's Dilemma

In 1950 an experiment was conducted by Merrill Flood and Melvin Dresher at the RAND corporation as an empirical test of the theoretical predictions of game theory. Two individuals played a version of the Prisoner's Dilemma with each other one hundred times. The payoffs (in pennies) were:

	B	1 (left)	2 (right)
A	1 (up)	-1, 2	1/2, 1
	2 (down)	0, 1/2	1, -1

with row choosers payoffs first and column choosers second. Each player has a dominant strategy. Row chooser is better off playing down no matter what column chooser does. Column chooser is better off playing left, regardless of row chooser's choice. If each player is rational, and plays dominant, the resulting payoffs are 0 and $\frac{1}{2}$ respectively. If each player made the other choice, up and right, each would be better off, row chooser getting $\frac{1}{2}$ instead of 0 and column chooser getting 1 instead of $\frac{1}{2}$. This is not what was seen in the experiment. The Nash equilibrium <down, left> is only played 14 times, while <up, right> is played 60. Flood concludes that game theory does not hold up as an empirical theory: "It seems unlikely that the Nash equilibrium point is in any realistic sense the correct solution ..." [Flood p.20]

He showed this report to John Nash, and Nash disagreed with his analysis. Flood reprinted Nash's comments in a long footnote:

"The flaw in this experiment as a test of equilibrium point theory is that the experiment really amounts to having the players play one large multistage game. One cannot just as well think of the thing as a sequence of independent games as one can in zero-sum cases. There is much too much interaction, which is obvious in the results of the experiment.

“Viewing it as a multistage game a strategy is a complete program of action, including reactions to what the other player has done. In this view it is still true that the real absolute equilibrium point is for A always to play 2, B always 1. [i.e Row Down, Column Left]

“However, the strategies:

A plays 1 ‘till B plays 1, then 2 ever after,

B plays 2 ‘till A plays 2, the 1 ever after,

are very nearly at equilibrium and in a game with an indeterminate stop point, or an infinite game with interest on utility, w it is an equilibrium point. ...”

Nash went on to argue that there is a case for applicability of the results to long finite repetitions:

“Since 100 trials are so long that the Hangman’s paradox cannot possibly be well reasoned through on it, it’s fairly clear that one should expect an approximation to this behavior which is most appropriate for indeterminate end games with a little flurry of aggressiveness at the end and perhaps a few sallies, to test the opponent’s mettle during the game.”

In this remarkable passage, Nash (1) invokes infinitely repeated games, (2) introduces the two interpretations of geometrical discounting in repeated games (either there is a constant probability less than 1 of one more play, “an indeterminate stop point”, or there is discounting of future payoffs in an infinitely repeated game, (3) gives a “grim trigger” strategy for maintaining cooperation, (4) introduces the idea of an epsilon-equilibrium, “almost an equilibrium”, and (5) invokes “bounded rationality” as a justification for applying this infinitely repeated game result to a fixed finite series of 100 trials. All of these ideas have enjoyed further development.

Nash’s explanation of cooperation in the Prisoner’s Dilemma is a “how possibly” explanation of co-operation in several ways. First, there are many possible equilibria in the repeated game, including “always defect”, the Nash equilibrium in the single one-shot Prisoner’s Dilemma. Second, the availability of a cooperative equilibrium in the infinitely repeated game depends on the rate at which the future is discounted. Players must put enough weight on the future, in order for cooperation to be sustained; the prospect of punishment must more than over-balance the prospect of immediate gain. Third, many punishment strategies other than grim trigger can sustain cooperation.

3 Infinite Repeated Games

There are two ways to get well-defined payoffs for strategies in infinitely repeated games. One is to take the route suggested by Nash, to geometrically discount the future. Then one can conveniently take the infinite sum. The other is to take an average. The payoff in an infinite sequence of outcomes is then the limit of the mean payoffs in initial segments. A number of folk theorems have been proved using the limit of the means definition.

Under both definitions, the model seems remote from human affairs, where repetitions are finite. But arguably the limit of the means definition is more remote, since the outcome of any finite initial sequence of repetitions, no matter how long, is irrelevant to the final value. For this reason, we confine the discussion here to Nash's version.

Consider how this works with the Prisoner's Dilemma. Grim Trigger strategies assure cooperation, by the threat of continued punishment for violation. Sufficiently high weight to the future will make punishment costly enough to outweigh any immediate benefit from defecting from co-operative play. Grim Trigger against Grim Trigger is a Nash equilibrium. Furthermore,

the threat of punishment is credible. If the trigger were to be pulled, carrying out the punishment would result in an equilibrium in the resulting subgame. That is, Grim Trigger against Grim Trigger is a subgame perfect Nash equilibrium.

Evidently this will work for a great variety of games - not just the Prisoner's dilemma. And a great variety of patterns of play can be supported by trigger strategies. Roughly speaking, you can get whatever you want in equilibrium in infinitely repeated games. This is the Folk Theorem that everyone knew, but no one bothered to write out for a while. It can be made precise in a number of ways.

Friedman (1971) published a folk theorem along these lines. The game to be repeated need not be 2-person; it can be an n-person game with a finite number of strategies. The point to be maintained as an equilibrium in the repeated game is any that is Pareto superior to a Nash equilibrium of the one-shot game. The threat is to revert to repeatedly playing that Nash equilibrium of the one-shot game. This is a trigger strategy with Nash equilibrium threats. The theorem says that if the discount rate is high enough - if enough weight is put on the future, such a point is an equilibrium in the repeated game. It is subgame perfect equilibrium. The threats are credible, because in the subgame that would ensue if the threats were triggered,

the players would be playing a Nash equilibrium of the ensuing subgame.

Friedman further generalizes the argument in a way quite material to the “Game of Life”. The stage games that are strung together to form the repeated game need not be iterations of the same game. The discount rates need not be the same from game to game. All that is required is that certain bounds be respected so that a large enough minimal discount gives a punishment that outweighs all temptations to defect. He also points out that games with different numbers of players can be covered by adding “dummy players” to those with a smaller number, to bring them all up to the same size. The dummy players get the same payoff no matter what they do, and have no effect on the payoffs of the real players.

Fudenberg and Maskin (1986) proved a stronger folk theorem for n-player repeated games. Instead of punishment by reverting to a Nash equilibrium, it relies on Minimax punishment. A player deviating from cooperative play is punished by the others by playing so as to minimize his payoff, no matter that he does. They show that any point that dominates the minimax point can be maintained as a subgame perfect equilibrium, provided the players place enough weight on the future. Subgame perfection with minimax threats does not come for free as with Nash equilibrium threats, and subgame perfection requires a more complex punishment strategy. Punishers are rewarded. There is also a technical

requirement (full dimension) for games of more than two players to make it possible to single out an individual defecting player for punishment.

The foregoing Folk Theorems all assume that all players observe everything that is done. If someone defects from cooperative play, everyone sees this and can react to it to impose punishment. This is an idealization that is evidently quite remote from the Game of Life, except in very special circumstances. There are, however, Folk Theorems for the case where there is only some noisy publicly observable signal of everyone's moves.

Fudenberg, Levine and Maskin (1994) prove two: one for Nash threats and one for minimax threats. Both, like the foregoing, provide strategies that form a subgame perfect equilibrium. There are requirements on the noise, more stringent for the minimax version.

In all of this, infinite horizon and arbitrarily high weight on the future play crucial roles in giving the general results. Yet humans, human societies and presumably humanity itself, have finite lifespans.

Finitely Repeated Games

Consider the other ideas that Nash put forward in his remarks to Flood.

He says that the trigger strategies "... are very nearly at equilibrium." This is made precise in Radner's (1986) concept of an *epsilon equilibrium*. An epsilon equilibrium (in a 2-person game) is a pair of strategies "such that each person's strategy is within epsilon of being a best response to the other's." The distance is measured in average payoffs. Although the best response to Grim Trigger in Flood's game - that is "Cooperate until your opponent defects and then defect forever - calls for a strategy that cooperates for 99 rounds and defects on the last, the increase in net average payoff resulting is small.

Here the difference is measured for strategies for the whole game. It may be objected that epsilon changes in subgames as the number of remaining repetitions gets smaller. Radner suggests an analogue of subgame perfection. This requires that the strategies get within epsilon of best response for the rest of the game for every continuation. He remarks that in this case there are epsilon equilibria in which players cooperate if there are sufficiently many periods remaining, and then stop cooperating close to the end.

How this generalizes is explored by Fudenberg and Levine (1986). They give general conditions under which a subgame perfect equilibrium in

an infinitely repeated game is a limit of epsilon-equilibria in finitely repeated approximations to that game.

Nash also invokes bounded reasoning capacity: ““Since 100 trials are so long that the Hangman’s paradox cannot possibly be well reasoned through on it ...”. This has been made precise by Neyman (1985) where games are played by finite automata. The strategies that are implementable are limited by the sizes of the automata. Some strategies in the finitely repeated Prisoner’s Dilemma, for instance Tit-for-Tat, can be implemented by a small automaton. If the number of rounds is large relative to the size of the automaton, strategies such as *Tit-for-Tat until the last round, then defect* cannot be implemented. The difference need not be very big for Tit-for-Tat to be a Nash equilibrium strategy. It is not subgame perfect, as the choice for the last round is just the choice in a one-shot Prisoner’s Dilemma. If the agents recompute at every stage of the finitely repeated game, then it is again possible to have initial cooperation and endgame defection.

This idea of analyzing repeated games using computationally restricted agents has been further developed in a large research literature. Notably, Binmore (1987,1988) considers play by automata that each run models of the others’ reasoning processes. This idea is still being explored. Critch (2019) analyzes a model in which artificial agents can read each others’ source code, with quite surprising results. Regarding the Folk Theorem,

Papadimitriou and Yannakakis (1994) produce a Folk theorem for finitely repeated games played by finite automata, where all payoffs that Pareto dominate the minimax point can be approximated arbitrarily closely.

A third approach moves from games of complete information to the games of incomplete information of John Harsanyi (1967-68). New possibilities emerge. Players do not necessarily know other players' payoffs, or know that others know theirs. All sorts of irrational behavior can be modeled, as Harsanyi suggests, as rational behavior of an actor with alternative payoffs. In this setting, given the right beliefs, we can even have cooperation in finitely repeated Prisoner's Dilemma. This was pointed out by Kreps and Wilson (1982) and by Milgrom and Roberts (1982) in two papers submitted at the same time to the same journal. Both were published with an introduction by all four in Kreps, Milgrom, Roberts and Wilson (1982). They point out that not only when each player thinks the other is irrational, but also with higher level failures of common knowledge, cooperation may be supported. For example, I may know that you are rational but want to trick you into thinking that I am not, so that you will cooperate in the short run. In all these scenarios, cooperation again breaks down in the endgame, provided that players know their own payoffs, but it is possible that cooperation can be sustained for long periods in finite repetitions of Prisoner's Dilemma.

The question of how Kreps, Wilson, Milgrom and Roberts generalizes was answered by Fudenberg and Maskin (1986). They develop a *Folk Theorem for games with incomplete information* in parallel with the one for infinitely repeated games with discounting. In this context, ruling out incredible threats is achieved by the requirement that the Bayes-Nash equilibrium is sequentially rational. The folk theorem is of the same strength. (Extension to three or more players requires the same technical condition, full dimension, as in the infinite case.) That is, any individually rational point – one that Pareto dominates the minimax point – can be approximated arbitrarily closely if the number of repetitions is great enough.

All in all, the prospects for interpreting the Game of Life as a bargaining game, via finitely repeated games, look promising. There are, no doubt, loose ends that need to be wrapped up. But there is a set of solid results that are in favor of the case.

5 Finitely Repeated Games, The Game of Life and the Game of Morals

The challenge of infinity for showing the Game of Life to be a grand bargaining game has been met in three ways for finitely repeated games: (1) bounded reasoning, (2) epsilon equilibria and (3) games of imperfect information. Each can be thought of as adding realism to the discussion.

They are not mutually exclusive, and have the potential for combining in perhaps synergistic ways.

This is not quite the end of the story. When we think of someone playing the “Game of Morals” – a social planner, or just an individual taking the moral point of view – we are thinking not just of selecting a point on the Pareto frontier of a Nash Bargaining Game. We are thinking of *how* it is sustained as such a point by our account of finitely repeated games. Our three approaches, taken individually, each put the question in a different light.

Consider approach (3) through games of imperfect information. What is being chosen in the Game of Morals? It is a set of prior degrees-of-belief. These may be beliefs that others are not rational, or that others do not think that that you are. Or that others have idiosyncratic values, or think you do, and so forth. This makes the Game of Morals seem like a way of picking an ideology. A social planner might be picking a propaganda campaign. This is very far in spirit from Binmore’s original intent. (A referee, however, remarks that it may not be so far from a Rawlsian “sense of justice” that supports and is supported by certain social institutions.)

On the other hand, (1), i.e. bounded reasoning, simply restricts the available strategies on the basis of cognitive complexity. Agents still can be

rational, and indeed have common knowledge of rationality. The Game of Morals still picks strategies that support the desired payoff profile as a Nash equilibrium of the repeated game. The bounds on reasoning can be taken empirically from the members of society, with some caveats if the members are taken to include corporations. Of course, finite does not mean small, and the finite number of repetitions to cover all cases can be arbitrarily large. But with that grain of salt, alternative (1) seems to fit Binmore's program fairly well.

Alternative (2), with epsilon best response and epsilon-equilibria, seems to fall somewhere in the middle. Legislating epsilon is dubious, although there could be a political-education or propaganda angle here. Epsilon could be taken to be descriptive, but in reality, it is clear that different individuals have different epsilons. This will have consequences that need to be explored. A social contract based on an epsilon that is too large would tend to be unstable.

As a combined approach, the priors used in (3) could be taken as objectively based on knowledge of typical reasoning power and typical epsilon among members of the population. Much remains to be explored in this direction. Prospects for precise results seem most promising for populations of interacting artificial agents.

References:

Aumann, R. (1997) "Rationality and Bounded Rationality" *Games and Economic Behavior* 21: 2-14.

Binmore, K. (1994) *Game Theory and the Social Contract vol. I: Playing Fair* MIT Press.

Binmore, K. (1998) *Game Theory and the Social Contract vol. I: Just Playing* MIT Press.

Binmore, K. (1987) "Modeling Rational Players I" *Economics and Philosophy* 3: 179-214.

Binmore, K. (1988) "Modeling Rational Players II" *Economics and Philosophy* 4: 9-55.

Critch, A. (2019) "A Parametric Resource-bounded Version of Löb's Theorem and a Robust Cooperation Criterion for Open-Source Game Theory" *Journal of Symbolic Logic* 84:1368-1381.

Flood, Merrill (1952) "Some Experimental Games" Working Paper RM-789-1
The Rand Corporation Santa Monica, California.

Friedman, J. (1971) "A Noncooperative Equilibrium for Supergames" *Review of Economic Studies* 38: 1-12.

Fudenberg, D. and Levine, D. (1986) "Limit Games and Limit Equilibria"
Journal of Economic Theory 38:261-279.

Fudenberg, D. and E. Maskin (1986) "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information" *Econometrica* 54: 533-554.

Fudenberg, D., Levine, D. and Maskin, E. (1994) "The Folk Theorem with Imperfect Public Information" *Econometrica* 62: 997-1039.

Harsanyi, J. (1967) "Games with incomplete information played by "Bayesian" players, Part I-III. Part I. The basic model." *Management Science* 14: 159-182.

Harsanyi, J. (1968a) "Games with incomplete information played by "Bayesian" players. Part II. Bayesian equilibrium points." *Management Science* 14: 320-334.

Harsanyi, J. (1968b) "Games with incomplete information played by "Bayesian" players. Part III. The basic probability distribution of the game" *Management Science* 14: 486-502.

Kreps, D., P. Milgrom, J. Roberts and R. Wilson (1982) "Rational Cooperation in the Finitely Repeated Prisoner's Dilemma" *Journal of Economic Theory* 27: 245-252.

Kreps, D. and R. Wilson (1982) "Reputation and Imperfect Information" *Journal of Economic Theory* 27: 243-279.

Milgrom, P. and J. Roberts (1982) "Predation, Reputation and Entry Deterrence" *Journal of Economic Theory* 27: 280-312.

Neyman, A. (1985) "Bounded Complexity Justifies Cooperation in the Finitely Repeated Prisoner's Dilemma" *Economics Letters* 19:227-229.

Papadimitriou, C. and Yannakakis, M. (1994) "On Complexity as Bounded Rationality" In [STOC '94: Proceedings of the twenty-sixth annual ACM symposium on Theory of Computing](#) 726-733.

Radner, R. (1986) "Can Bounded Rationality Resolve the Prisoner's Dilemma" In *Essays in Honor of Gerard Debreu* Ch 20 Elsevier 387-399.

Rubinstein, A. (1986) "Finite Automata Play the Repeated Prisoner's Dilemma" *Journal of Economic Theory* 39: 83-96.