

1. Introduction

This thesis advances game theory by formally analysing the implications of some of its most stringent assumptions. The approach followed here consists in examining the consequences of replacing some of the assumptions made in game theory for the sake of mathematical tractability with alternatives that –at least in some contexts– are more plausible. The method employed to conduct this research has been a symbiotic combination of computer simulation and mathematical analysis. Our results suggest that some of the most fundamental assumptions embedded in game theory may have deeper philosophical implications than commonly assumed.

1.1. Motivation

The value of advancing game theory seems clear: it is widely agreed that game theory has become one of the cornerstones of the social sciences (Hargreaves Heap and Varoufakis, 1995). There are widespread claims that it “provides solid microfoundations for the study of social structure and social change” (Elster, 1982), and that it “may be viewed as a sort of umbrella or ‘unified field’ theory for the rational side of social science” (Aumann and Hart, 1992). More recently, Gintis (2000) has stated that “game theory is a universal language for the unification of the behavioral sciences”. Even in the biological sciences it has been argued that some game theoretical concepts represent “one of the most important advances in evolutionary theory since Darwin” (Dawkins, 1989).

However, while extremely informative, game theory is at present somewhat limited in the sense that it is dominated by assumptions of full rationality, it generally ignores the dynamics of social processes, and it often requires disturbing and unrealistic hypotheses about individuals’ assumptions about other individuals’ cognitive capabilities and beliefs in order to derive specific predictions. Furthermore, it is often the case that even with heroic assumptions about the computational power and beliefs that every individual attributes to every other individual, game theory cannot reduce the set of expected outcomes significantly.

Thus, whilst acknowledging that the work conducted in game theory has been tremendously useful, a growing inter-disciplinary community of scientists think the time has come to extend game theory beyond the boundaries of full rationality, common-knowledge of rationality, consistently aligned beliefs, static equilibria, and long-term convergence. These concerns have led various researchers to develop formal models of social interactions within the framework of game theory, but relaxing its most stringent assumptions. Such models are providing not only valuable insights for the specific questions they address, but also the basis to assess how robust the results obtained in classical game theory are. This thesis is a contribution to this emergent programme of research.

1.2. Aim, approach and methodology

The overall aim of this thesis is to advance non-cooperative game theory by formally studying the implications of some of its assumptions that have been made for the sake of tractability and are not generally supported by empirical evidence. This has been done following two approaches:

- The first approach consists in examining the formal implications of replacing the unsupported assumptions in mainstream non-cooperative game theory relating to individual decision-making with assumptions that stem from empirical research. In particular, this thesis abandons the assumptions of complete information, common knowledge of rationality and consistently aligned beliefs, and contemplates instead members of two classes of decision making algorithms that have received strong support from cognitive science research: reinforcement learning and case-based reasoning.
- The second approach is used to extend mainstream evolutionary game theory. It consists in exploring the implications of a wide range of competing assumptions –all of them consistent with the essence of the theory of evolution– within a common framework. The results obtained using different assumptions are then contrasted in a coherent and systematic way.

In terms of methodology, there are four features that distinguish the work conducted in this thesis from most of the previous research undertaken in the same emerging field.

- First, the contributions made in this thesis have been placed in an overall framework that can encompass, in admittedly broad terms, most of the research conducted in game theory until now. This has permitted a more transparent comparison between the assumptions investigated here and those that have been addressed so far, and also between the results derived from this research and those obtained under other assumptions.
- Secondly, in terms of method, since most of the assumptions investigated in this thesis have not been formulated to allow for mathematical tractability, but to advance our formal understanding of social interactions in real life, new methodologies have had to be employed to supplement mathematical analyses. In particular, computer simulation has been used extensively to enhance and complement mathematical derivations. These two techniques have been combined in a way that is not common in the literature of game theory or in the field of social simulation. To be specific, most of the simulations reported in this thesis are just small advances at the edge of theoretical understanding. They are advances sufficiently small so that simplified versions of them (or certain aspects of their behaviour) can be fully understood in mathematical terms –thus retaining analytical rigour–, but they are steps large enough to significantly extend our understanding beyond what is achievable using the most advanced mathematical techniques available. In this way, simulations will be shown to extend theoretical knowledge in a rigorous, formal, and almost continuous way (Probst, 1999).
- The symbiotic use of mathematical analysis and computer simulation has allowed us to characterise both the short-term and the long-term dynamics of the models investigated in this thesis. This is in contrast with most game theoretical research –which is most often concerned with the identification of asymptotic equilibria– and with most research in the field of social simulation –which is often mainly concerned with the short-term dynamics.
- Finally, a great effort has been made to ensure that all models and simulations reported in this thesis can be easily scrutinised, used, replicated and reimplemented by independent researchers. In particular, all the computer programs used to conduct the research presented here have been released under the GNU general public licence (GPL), which is one of the licences

that scores best against the criteria set out by Polhill and Edmonds (2007) for releasing scientific software. GNU GPL grants the right to inspect, copy and distribute the source code, to modify it, and also to copy and distribute any modifications. It also guarantees that any modifications will be issued under a licence that preserves these rights (i.e. copyleft protection). Furthermore, following Polhill and Edmonds' (2007) guidelines, a substantial amount of work has been devoted in this thesis to *facilitate* the process of scientific critique of this research, by carefully commenting the code, providing extensive documentation, and creating several user guides for all the developed software. All the programs and documentation are included in the Supporting Material of this thesis.

1.3. Overall framework and specific contributions

To appreciate more precisely the specific contribution of this thesis to human knowledge, it becomes necessary to formalise some terms related to game theory first. In this thesis, a clear distinction is made between game theory used *as a framework*, and the different branches of non-cooperative game theory as we know them nowadays – e.g. classical game theory, evolutionary game theory and learning game theory.

Game theory as a framework is a methodology used to build models of real-world social interactions. The result of such a process of abstraction is a formal model that typically comprises the set of individuals who interact (called *players*), the different choices available to each of the individuals (called *strategies*), and a *payoff* function that assigns a (usually numerical) value to each individual for each possible combination of choices made by every individual. In most branches of game theory, payoffs represent the preferences of each individual over each possible outcome of the social interaction. The notable exception is evolutionary game theory, where payoffs most often (but not always) represent Darwinian fitness.

The feature of the social interaction to be modelled that makes game theory a particularly useful framework to employ is its *strategic* nature: the fact that the outcome of the interaction for any individual player generally depends not only on

her own choices, but also on the choices made by every other individual. Thus, game theory could arguably be defined as “the theory of interdependent decision-making” (Colman, 1995, pg. 3).

Game theory *used as a framework* provides a formal description of the social setting where the players are embedded. Importantly, it does not account for the players’ behaviour, neither in a normative nor in a positive sense. It is just not the realm of game theory *as a framework* to do so. It is only when different assumptions about how players behave –or should behave– are included in the framework, that game theory *as a framework* gives rise to the different branches that compose game theory as we know it nowadays. Here we outline the main features of the three most developed branches of deductive non-cooperative game theory at this time:

Classical game theory: Classical game theory was chronologically the first branch to be developed (Von Neumann and Morgenstern, 1944), the one where most of the work has been focused historically, and the one with the largest representation in most game theory textbooks and academic courses. Classical game theory is a branch of mathematics devoted to the study of how instrumentally rational players should behave in order to obtain the maximum possible payoff in a formal game.

The main problem in classical game theory is that, in general, rational behaviour for any one player remains undefined in the absence of strong assumptions about other players’ behaviour. Hence, in order to derive specific predictions about how rational players should behave, it is often necessary to make very stringent assumptions about everyone’s beliefs (e.g. common knowledge of rationality) and their interdependent consistency. Since such strong assumptions rarely hold in the real world, it is not surprising that when game theoretical solutions have been empirically tested, disparate anomalies have been found (see, for example, work reviewed by Colman (1995) in chapters 7 and 9, Roth (1995), Ledyard (1995), and Camerer (2003)). To make matters worse, even when the most stringent assumptions are in place, it is often the case that several possible outcomes are possible, and it is not clear which –if any– may be achieved, or the process through which this selection would happen. Thus, the general applicability of

classical game theory is limited. A related limitation of classical game theory is that it is an inherently static theory: it is mainly focused on the study of end-states and possible equilibria, paying hardly any attention to how such equilibria might be reached.

Evolutionary Game Theory: Some time after the emergence of classical game theory, biologists realised the potential of game theory as a framework to formally study adaptation and coevolution of biological populations (Lewontin, 1961; Hamilton, 1967). For those situations where the fitness of a phenotype is independent of its prevalence, optimisation theory is the proper mathematical tool. However, it is most common in nature that the fitness of a phenotype depends on the composition of the population (Nowak and Sigmund, 2004). In such cases game theory becomes the appropriate framework.

In evolutionary game theory, players are no longer taken to be rational. Instead, each player –most often meant to represent an individual animal– always selects the same (potentially mixed) strategy¹ –which represents its behavioural phenotype–, and payoffs are usually interpreted as Darwinian fitness. The emphasis is then placed on studying which behavioural phenotypes (i.e. strategies) are stable under some evolutionary dynamics, and how such evolutionary stable states are reached. Despite having its origin in biology, the basic ideas behind evolutionary game theory –that successful strategies tend to spread more than unsuccessful ones, and that fitness is frequency-dependent– have extended well beyond the biological realm.

The main shortcoming of mainstream evolutionary game theory is that it is founded on assumptions made to ensure that the resulting models are mathematically tractable. Most of the work assumes one single infinite and homogeneous population where players using one of a finite set of strategies are randomly matched to play an infinitely repeated 2-player symmetric game. In the last few years various alternative models (e.g. finite populations, stochastic

¹ This assumption, which is not always made in models of *cultural* evolution, is explained in detail in chapter 2.

strategies, multi-player games, structured populations) are being explored, but unsystematically.

Learning game theory: Like evolutionary game theory, learning game theory abandons the demanding assumptions of classical game theory on players' rationality and beliefs. However, unlike its evolutionary counterpart, learning game theory assumes that individual players adapt, learning over time about the game and the behaviour of others (e.g. through reinforcement, imitation, or belief updating). This learning process is *explicitly* modelled (Vega-Redondo, 2003, pg. 398). These investigations are being undertaken experimentally and formally (both analytically and using computer simulation), and special emphasis is being paid to the study of backward-looking learning algorithms, which seem to be more plausible than the forward-looking methods of reasoning employed in classical game theory. The latter appear to be very demanding for human agents (let alone other animals) and remain undefined in the absence of strong assumptions about other players' behaviour and beliefs. Some of the most studied classes of decision-making algorithms in the literature are: reinforcement learning (with experimental studies conducted by e.g. Erev et al. (1999), theoretical work done by e.g. Bendor et al. (2001b), and studies of the dynamics carried out by e.g. Macy and Flache (2002)), belief learning (with theoretical work on fictitious play developed by e.g. Fudenberg and Levine (1998)), and the EWA (Experience Weighted Attraction) model (Camerer and Ho, 1999), which is a hybrid of reinforcement and belief learning.

This thesis makes two specific contributions to the development of learning game theory and one in the field of evolutionary game theory. The first contribution to learning game theory is to elucidate the implications of assuming that players use a simple form of reinforcement learning as decision-making algorithm. Reinforcement learning, being one of the most widespread adaptation mechanisms in nature, has attracted the attention of many scientists and engineers for decades. This interest has led to the formulation of various different models and –when feasible– to the theoretical analysis of their dynamics. This thesis provides an in-depth analysis of the transient and asymptotic dynamics of one of the best known

stochastic models of reinforcement learning (Bush and Mosteller, 1955) for 2-player 2-strategy games.

The second contribution to learning game theory is a detailed exploration of the implications of case-based reasoning as decision-making approach in strategic contexts. Case-based reasoning consists in “solving a problem by remembering a previous similar situation and by reusing information and knowledge of that situation” (Aamodt and Plaza, 1994). Case-based reasoners do not employ abstract rules as the basis to make their decisions, but instead they use similar experiences they have lived in the past. Such experiences are stored in the form of cases. The distinguishing feature of case-based reasoning as problem-solving mechanism is that “thought and action in a given situation are guided by a single distinctive prior case” (Loui, 1999). To our knowledge, the implications of this type of reasoning in strategic contexts have not been explored before.

Finally, the contribution of this thesis to evolutionary game theory is a systematic exploration of the impact of various assumptions made in this field; this exploration is undertaken by studying the structural robustness of evolutionary models of cooperation using a computational tool built for this specific purpose: EVO-2x2. EVO-2x2 is a computer simulation modelling framework designed to formally investigate the evolution of strategies in 2-player 2-strategy (2x2) symmetric games under various competing assumptions.

A significant part of the work conducted in this thesis is sufficiently general to be valid in a wide range of social interactions, but some of it has had to be focused on particular types of social interactions. Whenever there has been a need to select a specific type of social interaction to investigate (even if the only purpose was to illustrate the applicability of more general findings), we have always studied social dilemmas (Dawes, 1980). Social dilemmas are social interactions where individual rationality leads to outcomes for which there is at least one feasible alternative preferred by everyone. In such situations, decisions that make sense to each individual can aggregate into outcomes in which everyone suffers (Macy and Flache, 2002). The focus of this thesis has been on social dilemmas because of their importance in the social and biological sciences, and because the predictions

of classical game theory in this context clash with widely shared intuitions and empirical results (see, for instance, work reviewed by Gotts et al. (2003b) and by Colman (1995) in chapters 7 and 9).

1.4. Outline of the thesis

The structure of this thesis is as follows: chapter 2 outlines the main assumptions made in game theory. We analyse each of the following branches in turn: game theory used as a framework, classical game theory, evolutionary game theory, and learning game theory. This critical review of the main assumptions made in deductive game theory will serve as a framework to clearly identify those assumptions that will be abandoned in the subsequent chapters of this thesis, and those that will be retained. Chapter 3 clarifies the scope of this thesis within game theory and explains social dilemma games in detail. It also describes the methods that have been used to formally analyse the models developed in chapters 4, 5 and 6. Chapter 4 is an in-depth analysis of the transient and asymptotic dynamics of the Bush-Mosteller reinforcement learning algorithm for 2-player 2-strategy games. Chapter 5 is an exploration of case-based reasoning as decision-making algorithm in strategic contexts. Chapter 6 describes EVO-2x2, the modelling framework developed in this thesis to assess the impact of various assumptions made in mainstream evolutionary game theory for the sake of mathematical tractability. The use of EVO-2x2 is illustrated by conducting an investigation on the structural robustness of evolutionary models of cooperation. Chapter 7 is a general discussion of the results obtained in chapters 4, 5 and 6. We also discuss the value of the models developed in this thesis, and how they could be validated. Chapter 8 summarises the main conclusions of this work and identifies areas for further research. The proofs of the theoretical results derived in this thesis can be found in the appendices. This thesis also comprises extensive supporting material, including the source code of every computer program we have used in this research, together with user guides and instructions to replicate every experiment reported here.