# Exploring Individual Preferences in Economic Contexts: Three Essays in Evolutionary Game Theory

Vorgelegt von

Niko Noeske

aus Bonn

Bielefeld 2011

Dekan: Prof. Dr. Herbert Dawid

Erstgutachter: Prof. Dr. Frank Riedel

Zweitgutachter: Prof. Dr. Herbert Dawid

i

# Acknowledgments

There are several people who I wish to thank for their plenary support and professional help. At the first place, I deeply thank Professor Dr. Frank Riedel for the opportunity he gave me to work as a doctoral student and for the continuous guidance and support he provided during the years. Whenever I asked him for help, he took his time to discuss my research agenda as well as my ambitions and plans. In addition, I thank Professor Dr. Herbert Dawid for his effort taken in surveying this thesis. Furthermore, I would like to thank the many seminar participants and colleagues in Bonn and Bielefeld for helpful comments, suggestions, and fruitful discussions. Special thanks, in this regard, go to Daniel Wiesen, Lars Metzger, and Marcelo Cadena.

# Contents

# Chapter 1

# General Introduction

My thesis is divided into four parts. I start with surveying basic facets and differences of orthodox and evolutionary game theory which provide the methodological tools of use. After this, I clarify the meaning of the concept of "individual preferences" and explain the indirect evolutionary approach and its specific use in this thesis. The remaining three chapters comprise three self-contained essays yet connected by the same topic of indirect evolution of individual preferences in economic contexts.

## 1.1 Orthodox and Evolutionary Game Theory

Game theory is the mathematical analysis of interdependent decisions and outcomes, in the sense that those who make decisions are affected by their own choices and by the choices of others. In terms of game theory, a *game* is a strategic situation in which two or more *players* interact with each other by choosing *strategies* from a bundle of alternatives which calculate their *payoffs* or *utilities*. Accordingly, a game can completely be represented by the triple $G = (\text{players}, \text{ strategies}, \text{ payoffs})$. Typically, such a triple represents a game in normal (or strategic) form where the players

choose their strategies simultaneously.[1]

The interpretation of a game may differ with respect to peculiarities of the model under consideration. For example, in one setting, the underlying model describes an industry structure where players are firms which compete with each other by adjusting their decisions on the amount of output they will produce—or, firms adjust their decisions on the prices of the goods they sell. In a different class of games, one might think of players as fishermen who exploit an inshore fishery resource by adjusting their fishing efforts. In the case that players put too much effort in the game, the stock of fish is reduced to a level where it can hardly be recovered. But in the case that players put too little effort in the game, they relatively fail by losing profit comparing with their fisher colleagues. So, the difficulty of the fishermen is to adjust their inputs by reasoning on the fisher colleague's inputs and the recovering rate of the fish stock. Of course, the situation of the fishermen is basically a metaphor and can easily be transferred to the fundamental problems of modern times like the environmental disasters or financial crises. This thesis explicitly and implicitly attends to the mentioned (and other) sorts of problems by assuming interdependent success in the manner of the material payoffs as given in the first two essays.

In a different setting, the players are, for example, politicians who invest high efforts or money to get elected at the next election. The particular investment is the result of an evaluation process which integrates both the winning probability as an increasing function of own input and the loss of investment. In a comparable strategic setting, one might think of companies which invest in research and development (R&D) since they may benefit from introducing a new product first by earning some monopoly rent before a competitor enters the market with a similar product. These strategic considerations describe just two specific examples of a broad class of

---

[1]My thesis focuses on simultaneous decisions. If one wants to explicitly account for a time (or sequential) structure within the decision period, i.e. player $A$ knows the decision of player $B$ before player $A$ makes his decision, one should model games in extensive form.

games known as "contests" where players make irrecoverable investments in order to influence the probability of winning a certain price. The last essay of this thesis deals with such strategic settings.

A standard analytical procedure of orthodox game theory (henceforth OGT) is the usage of a Nash equilibrium to examine players' strategic choices in social life. A Nash equilibrium describes a profile of chosen strategies where no player benefits from unilaterally changing the chosen strategy since it is a best response given the other players' strategies.[2] However, a Nash equilibrium does not necessarily imply that corresponding payoffs are optimal. For example, by standard assumptions of the well-known prisoners' dilemma, where players can behave either cooperatively or non-cooperatively, "defection" (behaving non-cooperatively) is a Nash equilibrium but yields suboptimal payoffs. In games with a unique Nash equilibrium, one can rightly argue that Nash's concept is a powerful analytical tool to calculate the actions of a game and subsequent outcomes if the players are full-fledged individuals who know all the details of the game and their opponents and use these details to assess their decisions. However, as the previous sentence suggests, there are two apparent shortcomings of the Nash equilibrium concept.

First, there are sometimes many equilibria in a game that accord to the characteristics of Nash's concept. The emergence of more than one equilibrium raises the question of which one will actually arise in a certain conflict. The existence of more than one Nash equilibrium in games has led research in non-cooperative game theory to develop refinement concepts which give the equilibria a stronger justification for evaluating players' strategies and subsequent outcomes (most exemplary is Reinhard Selten's concept of "trembling hand" where players play "off-the-equilibrium" with a small probability). The attempts which aim at finding Nash's refinements have guided to a vast number of concepts which justifies nearly any Nash equilibrium by

---

[2]By extending John von Neumann and Oskar Morgenstern's (1944) pioneering work on game theory, John Forbes Nash, Jr., (1950) conceived the notion of the equilibrium concept which bears his name and which has revolutionized game theory and economics.

the one or other interpretation of the particular game standards like the cognitive abilities of the players.

Second, standard requirements of the Nash equilibrium concept and many of its refinements include a perfect, common-knowledge rationality of the players. More precisely, the players are rational in the sense that they take all available information into consideration and choose actions that maximize their expected payoffs given that the other players are informed and act in the same way (which makes common-knowledge rationality implies that all players are perfectly rational in the same sense which subsumes, inter alia, that they know this fact). This is obviously a very strong postulate which has provoked many researchers to reconsider the usefulness of this assumption and to think of alternative approaches.

As a sidestep, experimental research has shown that humans do not act according to predictions of the strong assumptions on rationality, but rather to simple decision rules which are, for example, inferences from learning processes in real-life situations.[3] Apart from game theory, the existence of "limited agents" is not new in economic research. For example, Gérard Debreu (1959, p. 37) implicitely accounts for "imperfect" agents by writing in his fundamental "Theory of Value": "an agent is characterized by the limitations on his choice, and by his choice criterion". In other words, it is the individual opportunities that count. As soon as the experimental research programs manifest these findings, the rather static solutions of OGT obtained by calculating the behavior of perfect rational individuals have appeared to be fairly unrealistic. With these insights, equilibria are no longer necessarily considered as designated profiles in one-shot games which appear from the synthetic postulate of rational agents, but rather as the result of bounded rational players who play the same strategic situation over and over. This modified economic thought is represented by a relatively new branch of economic research which has its roots in

---

[3]Cf. Camerer (2003) for an overview. Gigerenzer and Selten (2001) use the metaphor of an "adaptive toolbox" to explain that decision makers are equipped with a certain bundle of strategies and that adaptive decision making arises from choosing from this particular bundle.

theoretical biology and the Darwinian 'survival of the fittest' doctrine.

Starting with John Maynard Smith and Georg R. Price's solution concept of an evolutionarily stable strategy (ESS)[4] in 1973 and, more rigorously, with the publication of John Maynard Smith's book "Evolution and the Theory of Games" in 1982, evolutionary game theory (henceforth EGT) has effected to attract the attention of many economists who doubt the classical concept of rational agents to examine human behavior in different strategic settings. Although EGT has its roots in biology, it is not defective that EGT has become of intensified interest to economics and social sciences in general.[5] This is because 'evolution' as processed by EGT is not necessarily biological evolution in the sense of gene transmission. Instead, 'evolution' may rely on cultural processes where values like conventions, norms, beliefs, or ideologies are shaped over rather short lengths of time.

The recent development of EGT is not only substantiated by biology but also by traditional economic game theory though. Foremost, it is by now well understood that Nash already had a dynamical population model in mind when developing his equilibrium concept. Likewise, Reinhard Selten's trembling hand concept opened the door for discovering the realm of bounded rational decision makers. This indicates that EGT is clearly a continuation of OGT which highlights the dynamical aspects of equilibrium selection and declines the necessity of "hyperrationality". Moreover, the fusion of evolutionary theory and economic theory has not come as a surprise since both fields use game theory as an analytical tool to sketch interdependence extensively. In fact, interestingly and somewhat ironically, as traditional game theory appears maybe more suited in biology where the players are either species or genes,[6] the recent approaches of EGT appear to be approximately suited to the field

---

[4]Roughly, an ESS is an incumbent strategy in a society ("population") that cannot be replaced by a rare mutant strategy under the influence of evolutionary pressure.

[5]Beside economics, EGT plays an increasing role in psychology, anthropology, sociology, as well as in philosophy to name the main areas only.

[6]In the preface of Maynard-Smith (1982), he writes: "Paradoxically, it has turned out that game theory is more readily applied to biology than to the field of economic behavior for which it was originally designed."

of behavioral economics as for evolutionary theory in the literal sense. Primarily, this is because economic agents are humans (or, e.g., companies which are driven by humans) who (which) are not as perfect as OGT suggests. Precisely, the advantage of EGT is that it is modest in comparison with the rationality postulate of OGT since it assumes agents to optimize behavior on a dynamic route of trial and error processes without requiring the agents to be rational at all.

## 1.2 Individual Preferences

Traditional economic science has built upon assumptions of *homo economicus*, or economic man, a rational and purely self-interested economic subject. However, there is a growing body of evidence in different research disciplines such as psychology, sociology, or economics that reveal the limitation of models with such assumptions. People often fail to maximize their own economic objectives because human choices are not only driven by material self-interest but also by "moral sentiments" or social norms like sympathy, compassion, guilt, or reciprocity. Many experimental studies in economics reveal that decision makers do not act as predicted by models of homo economicus. For example, consider the ultimatum game. The ultimatum game is one of the most intensively studied research subjects in behavioral economics. This is because it is simple and it is immediate to gain some insight into the economic psychology of behavior. It works as follows. There are two parties who have to agree on a fixed amount of money or other goods, say 100 units. The first mover has to make an offer $\delta \in [0, 100]$ to the second mover who either accepts or rejects the offer. In the case of accepting the proposer gains $(100 - \delta)$ and the responder gets $\delta$, while in the case of rejecting the game ends with zero payoffs for both parties. A relatively robust result throughout the experimental studies is that if $\delta$ is less than 30 (the proposer offers less than 30% of the pie), rejection is the usual consequence (cf. Güth et al., 1982; Camerer and Thaler, 1995; Roth, 1995,

and references therein). In addition, a large fraction of players make "fair" offers where both players get approximately the same. However, this is quite contrary to what orthodox equilibrium theory predicts: the unique subgame perfect Nash equilibrium is given by $(\delta = 0, \text{accept})$. A similar situation, which is even simpler than the ultimatum game, is known as the dictator game. A dictator divides an amount of money or other goods between himself and the receiver. Robust results can be summarized as follows. About 80% of the offers (the part which the receiver gets) are between zero and half of the pie. Roughly 20% offer zero, and offers larger than half of the pie are nearly never observed. Again, the usual behavior is far from "economic man" who does not make gifts without return (cf. Hoffman et al., 1996; Eckel and Grossman, 1998).

The crucial question is then what are the engines driving such behavior? Understanding the motivations of people is important because they determine the course of action which creates our social and economic world. This is where the conception of individual preferences is taken into consideration. In many studies, individual preferences are equivalent to social (or other-regarding) preferences where economic actors are concerned with the well-being of others as well as their own well-being (cf. Sobel, 2005, and the references therein for a recent development). In other studies, individual preferences include motivations that depart from models of homo economicus in that people have feelings of overconfidence, seek a high reputation, follow different ideologies, or are "economic irrational" for other reasons. A main field of individual preferences in the present thesis deals with social preferences where individuals take others' payoffs into consideration. In this realm, people are said to have "interdependent preferences" if they care about others' payoffs. An individual has positively (negatively) interdependent preferences if the payoff of others' positively (negatively) enters his utility. In contrast, an individual has independent preferences if his utility does not depend on others' payoff. In models comprising this kind of interdependence to explain certain behavior, it is usual to refer to positively (negatively) interdependent preferences as altruism (envy) and

7

to independent preferences as egoism. Since these notions play a central part in the following chapters, and there exist no definitions that are generally accepted in all scientific fields where interdependent (social) preferences appear,[7] it is necessary to explain how they should be understood in the present thesis. Firstly, if I write about social preferences, I will do so on the fundament of a trivial truism, namely that all people, whether endowed with altruistic, egoistic, or envious feelings, want to satisfy their own desire. For example, the altruist wants to help others and the envious person is maybe willing to harm others just to satisfy his feeling of not being envious anymore. If they engage in a game and if they have more than just one choice, both types will surely choose those actions that will maximize their own well-being (at the best of their knowledge). From this point of view, one can think of the whole world only consisting of fundamental egoists. However, in this thesis, it is a different level where the "evolution of preferences" happens. The fact that it is me who wants something is banal and not really worthwhile to think about; the essential point is what is it what I want. In other words, it is the content of our desires which defines what type of person we are. Secondly, social preferences are "intrinsic" motivated values which are subjective in the sense that only the person who owns that value is in the position to judge it. For example, an altruist seeks merely for a "warm glow" whose only utility he achieves stems from the act of giving (cf. Andreoni, 1990). An envious person does not explicitly strive for higher own payoffs but only for a subjective better feeling he gains from an improved relative standing in society. In contrast, "extrinsic" motivated values are objective in the sense of intersubjective measurable (e.g. the players strive for more money in economic contexts). A third point is that a person is not only altruistic, selfish, or envious but features all preferences. The specific circumstances of a strategic situ-

---

[7]For example, in evolutionary sociology and biology it is partly agreed to think about altruism as a trait which is essentially reciprocal in the sense that the evolutionary players anticipate (consciously or unconsciously) to gain from their "goodwill" in certain payoffs they expect from future conflicts, cf. Trivers (1971). In other studies it is common to think about altruism as a trait which does not claim for future benefits in the reciprocal sense, cf. Fehr and Fischbacher (2002).

ation (a game) develop the individual preferences' which determine the behavioral consequences and subsequent outcomes. In fact, addressing this issue will be an essential point in this thesis.

Beside the dimension of altruism (where envy is the negative part) there are two other dimensions of individual preferences that appear in this thesis: strong reciprocity and self-confidence. Reciprocity describes people's tendency to reward perceived kindness and to punish perceived unkindness. The adjective "strong" expresses the intrinsic value of reciprocity. In particular, strong reciprocity refers to the conception that people act reciprocal in order to satisfy their subjective fairness emotions and not to expect higher economic revenues in the future (cf. Gintis, 2000). The dimension of self-confidence varies from strong underconfidence to strong overconfidence. People are underconfident (overconfident) when they undervalue (overvalue) their own ability in certain situations (cf. Ando, 2004). The foregoing sentences describe how these psychological concepts should be understood in this thesis, however, the exact meanings of these concepts are both fully described and only palpable with the mathematical precision given by the well-being functionals of the three essays.

## 1.3 Indirect Evolutionary Approach

The indirect evolutionary approach is the central tool for exploring individual preferences in economic contexts in the thesis at hand. The methodology is applied in all essays of this thesis and further illustrated there—however, by virtue of its relevance, I give a brief introduction in this section. The indirect evolutionary approach combines traditional assumptions of economic decision making, as in OGT, with evolutionary processes resulting from EGT. In particular, it states that people behave rationally according to the subjective preferences they own but the success of these preferences are measured by evolution. The immense benefit of this approach

9

is that it allows the researcher to endogenize certain preferences in different contexts instead of just assuming them as given. Hence, indirect evolution seems to be more sophisticated than the usual direct approach (where a type is pre-programmed to play a certain strategy)[8] since subjective motivations initiate actions, and hence, constitute the more profound basic of economic decision making. There are some precursors of this approach (e.g. Becker, 1976), but subsequent research on individual preferences using indirect evolution stems from the methodological frame and pioneering work of Werner Güth and Menahem Yaari and their analysis on reciprocity in ultimatum bargaining in 1992. Seminal studies which use this approach are, among others, Bester and Güth (1998); Possajennikov (2000); Heifetz et al. (2007a); Leininger (2009).

There seems to be some agreement about the exact way of using the approach of indirect evolution which is considered in most of the corresponding studies: first, the individuals play a Nash equilibrium according to their subjective preferences. Second, the particular Nash equilibrium play is then inserted into an objective payoff function which measures the fitness of the subjective preferences of the players. Third, the ESS conception is then employed in order to derive stable biases from the created objective function. On the one hand it is comprehensible that the most famos solution concept of EGT is used to rationalize certain "economic irrational" behavior by certain preferences but on the other hand there are some apparent shortcomings of ESS in some cases which seem to be somewhat unvalued in the literature of indirect evolution. Foremost, ESS is a static concept which helps explain whether a somehow reached population state is immune to rare mutations but says only little about the evolving to such a state. In addition, Oechssler and Riedel (2001, 2002) show that ESS is an insufficient concept for the purpose of studying dynamic stability (which referes to the question whether any sufficiently small change of a population is such that the new population stays close and/or converges to the former population) of the replicator dynamics (where the reproductive success of a

---

[8]Cf. Weibull (1995) for a standard textbook treatment.

certain strategy is measured by the difference between the payoff of that strategy and the population average payoff) if the underlying strategy space is continuous. Since it is natural to assume that people perceive individual preferences from a continuum and that the biological replicator dynamics shape these preferences, it is worth considering alternative evolutionary concepts. In respect thereof, the first two essays of this thesis deal with solution concepts, in particular refinements of ESS, which appears more suited to the field of preference evolution. The last essay uses ESS for finite populations on the preference level—a concept from EGT which should not necessarily be seen as a refinement of Nash's equilibrium concept (cf. Schaffer, 1988).

# Chapter 2

# Altruism and Envy Revisited: On Evolutionary Stability with a One-dimensional Continuum

*Bester and Güth [Journal of Economic Behavior and Organization **34**, 1998, 193-209] take advantage of the "indirect evolutionary approach" and the evolutionary solution concept of an Evolutionarily Stable Strategy/State (ESS) to explain the evolutionary causality of altruism. However, ESS says only little about the evolving to such a state and becomes insufficient for dynamic stability with respect to the replicator dynamics if the underlying strategy space is continuous. We build on this work by allowing envy and adopting alternative solution concepts, namely Continuously Stable Strategy (CSS), Neighborhood Invader Strategy (NIS), and Evolutionary Robustness (ER). These concepts are much in line with the one-dimensional, continuous frame of opposed preferences competing with respect to the topology of weak convergence. The evolutionary qualities of altruism and envy are determined by the strategic environment. Moreover, we introduce an alternative definition of altruism and envy and show that the existence of a sophisticated perception of co-players' well-being is immediately negligible with respect to the evolutionary fitness of that preference.*

## 2.1 Introduction

A traditional assumption in economics is that agents are profit maximizers (egoists). However, this view is at odds with real life evidence and results of numerous economic experiments in several social settings. The behavior that individuals often show rather reflect so-called "other-regarding preferences"; for a comprehensive overview of this subject, see Fehr and Fischbacher (2003), and also Sobel (2005), along with the references therein.[1] Those findings that are so contradictory to the traditional matter of a homo economicus, whose adherent preference structure is purely determined by economic profit, have induced many researchers to consider altruistic and envious preferences, henceforth treated as endogenous features of a game.[2] An expedient technique in evolutionary game theory for exploring the evolutionary fitness of such intrinsic motivated values is given by the approach of indirect evolution, initiated by Güth and Yaari (1992). Given indirect evolution, the players maximize their perceived payoffs, i.e. their subjective preferences, with a strategic behavior that determines the objective game payoffs which, in turn, represent the evolutionary fitnesses. In this vein, one is able to draw conclusions for the evolutionary fitnesses of the underlying preferences in an indirect way.

Bester and Güth (1998) take advantage of this method to study the evolutionary fitness of altruism in symmetric 2-person games. In their work, altruism is identified by subjective preference functions that formulate the true well-being of the players by an individually weighted convex combination of own and opponents economic

---

[1]The overall aim of an egoist is supposed to be profit maximization in an economic sense (e.g. more monetary profit is better). Of course, one can think of an *egoist by definition* whose motivation is independent of the underlying utility since any motivations are always *personal* in the end. But this approach would only lead into trouble with the terminology of egoism (e.g. one has to think of altruistic egoists or egoistic egoists) and would apparently effect no extra gain for the analysis of the "evolution of preferences".

[2]For a wider discussion of envy in economics see Mui (1995) and the references therein. The study of altruism in economics is much more present in existing literature; see, e.g., the cardinal work of Becker (1976). For further models that endogenize certain preferences like fairness, status-concern, overconfidence, reciprocity, and morality, see, among others, Huck and Oechssler (1999), Fershtman and Weiss (1998), Kyle and Wang (1997), Guttman (2000), and Güth and Ockenfels (2005).

profit. They show that a key relation is given by the influence of altruism on the equilibrium actions of the co-players. Bester and Güth call this implication the *strategic effect of altruism.* The basic prerequisite to the validity of this effect is that preferences are (at least partially) observable, see also Güth and Peleg (2001), Ely and Yilankaya (2001), Ok and Vega-Redondo (2001), Heifetz et al. (2007b), and Dekel et al. (2007) for the necessity of common knowledge of preferences. The basic findings in these papers are that other-regarding preferences are evolutionarily viable in perfect common knowledge games and intermediate cases (e.g. the players know the distribution of preferences in the population so that they can expect a preference with a certain probability) but not in private information games. Following Bester and Güth, we will stick to the assumption that preferences are common knowledge in a perfect sense. Then, depending on the specific type of strategic game externality, the strategic effect points out that possessing altruistic preferences can either be harmful or profitable.[3]

The central issue of Bester and Güth enters into the question of the required circumstances for the evolutionary causality of altruism. Put differently: Why does altruism exist from an evolutionary point of view? The finding is that altruism is *evolutionarily stable* in the sense of an *evolutionarily stable strategy* (ESS) (Maynard-Smith and Price, 1973), i.e. an incumbent strategy that cannot be replaced by a rare mutant strategy under the influence of evolutionary pressure, if the game exhibits *strategic complements.* Otherwise, if the game exhibits *strategic substitutes*, only the single egoistic preference is ESS.[4]

Possajennikov (2000) and Bolle (2000) adopt the basic model of Bester and Güth but enlarge the space of available preferences by allowing the opposite of altruism,

---

[3]Of course, the strategic effect is relevant to *any* other-regarding preferences and not confined to altruism. Heifetz et al. (2007b) show very broadly that the strategic (or indirect) effect of other-regarding preferences, that influence the behavior of others in a certain manner, is generally quite stronger than the direct effect that naturally reduces one's own profit.

[4]Beside the basic finding, Bester and Güth identify two further interesting facts. Firstly, a population full of altruists is more successful than a population full of egoists; and secondly, concerning a game played by two different individuals, the less altruistic person is more successful.

which is referred to as "spite" in Possajennikov and "envy" in Bolle.[5] Both authors imply the ESS-concept as well to examine the evolutionary viability of those preferences. With the enlargement of the preference space, the egoistic preference does no longer comply with the requirements of ESS; instead, in the case of strategic substitutes, an envious preference does.

One major reason why the concept of ESS has reached such popularity is due to the fact that the requirements of ESS are sufficient for asymptotic stability with respect to the well-known replicator dynamics of Taylor and Jonker (1978), where players choose from a finite number of pure strategies. However, Oechssler and Riedel (2001) show that this fact is no longer necessarily true if the available strategy space becomes continuous. Even the stronger requirements of a strict Nash equilibrium are no longer sufficient to guarantee dynamic stability with respect to the replicator dynamics. Thus, following the probably most natural assumptions that preferences are available from a compact continuum and that a final population composition is unconsciously reached through evolutionary selection (as in the replicator dynamics) and not through rational decision making guides to the striking fact that stronger stability concepts than ESS are needed.

If one deals with dynamic stability and/or convergence in infinite strategy games, one has to think about the choice of the appropriate topology on the space of probability measures (which are the populations). The key question in this respect is: When is a population *evolutionarily close* to another population? In the discrete case, the choice of the appropriate topology is out of the question since two populations are always in Euclidean distance of each other. However, there are different alternatives of closeness in the continuous case. Two major options are framed by the *strong topology* (which is equal to the *variational norm*) on the one hand, and by

---

[5]We will refer to "envy" in this respect, partly for the sake of convenience. Also, we refrain from a broad psychological or philosophical debate about the opposite of altruism because we believe that the constraints of the model at hand allow for a wider range of terminologies. Thus, whenever the term "envy" arises, the field-equivalent is "spite" in Possajennikov (2000) and remains "envy" in Bolle (2000).

the *weak topology* on the other hand. There are several arguments for both choices.[6] The present paper is mainly engaged in the weak topology. Primarily, because the weak topology respects the natural resemblance of slightly different preferences.[7] In our setting, closeness is then certain in both evolutionary incidents: A large preference-shift by a sufficiently small fraction of the population and a sufficiently small preference-shift by a large fraction of the population.

Whilst the former happening is much in line with ESS the latter is more consistent with the attributes of a *Continuously Stable Strategy* (CSS, introduced by Eshel and Motro, 1981; Eshel, 1983) and a *Neighborhood Invader Strategy* (NIS, introduced by Apaloo, 1997). A CSS includes ESS and states that a sufficiently small homogeneous population-change from the ESS is such that strategies closer to the ESS are fitter than the new population strategy. A NIS is a strategy that is able to invade any sufficiently close, homogeneous neighbor via a higher fitness. From this perspective, the strong concept of *Evolutionary Robustness* ($\mathscr{ER}$, introduced by Oechssler and Riedel, 2002) unifies these approaches in the weak topology. A population is $\mathscr{ER}$ if it gains a higher than average payoff against all possible populations that are close in the weak topology. Our basic intention is to apply these concepts to the Bester/Güth game (or rather to an extended version in the sense of Bolle and Possajennikov). We find that the equilibrium preference satisfies the supplementary requirements.

The paper is organized as follows. The next section deals with the relevance of the finite and continuous strategy space concerning the replicator dynamics, the issue of a suitable topology, and gives a detailed presentation of the different stability concepts under study. In section 2.3, we review how indirect evolution of preferences works. Section 2.4 then introduces the basic evolutionary game and applies it to the stability concepts. Section 2.5 deals with what we call "sophisticated-perceptive

---

[6]For a detailed presentation and discussion, see Oechssler and Riedel (2001, 2002).

[7]This is because the corresponding *distribution functions* of the populations are then close to each other. See section 2.2 for details.

preferences" which extend other-regarding preferences in the sense that the players' take the real well-being, i.e. the subjective preferences, of others into account and not just the other players' economic profits. The main finding here is that the existence of such preferences is immediately negligible regarding the indirect evolutionary analysis. Finally, conclusions are drawn in section 2.6.

## 2.2   The Evolutionary Conception

The dynamics of evolutionary systems where the players choose their strategies from a continuous set are widely studied in the adaptive (or strategy) dynamics approach which follows the simplifying assumption that each population is homogeneous in the strategy choice and remains so during the course of evolution (cf. Marrow et al., 1996; Abrams, 2001). If one drops this simplifying assumption and considers a dynamical system where the aggregate play of a population is described by a distribution on the strategy space then one reaches (e.g.) the replicator dynamics with a continuous strategy space. It is the latter approach and the relevance of CSS, NIS, and $\mathscr{ER}$ which we mainly focus on; however, these concepts also give some insight regarding the adaptive dynamics approach. But first, we review the ESS concept and the replicator dynamics in the standard finite strategy space.

### 2.2.1   Finite Strategy Space

The most popular solution concept used in evolutionary game theory is known as ESS, which has been introduced by the biologists Maynard-Smith and Price (1973) for 2-person normal form games. To examine ESS, suppose that *evolutionary fitness* (or *reproductive success*) is defined by an individual's payoff resulting from evolutionary agents repeatedly drawn at random from one large population competing in pairwise contests. If all members of this population play an ESS, say $x^*$, then a small injection of mutants exerting any (pure or mixed) deviant strategy $x \neq x^*$

17

is initially less successful than the incumbent strategy $x^*$ and will eventually disappear. The exact impact of ESS, for *evolutionarily stable strategy*, is given by the following two well-known conditions.

$x^*$ is ESS iff, for all $x \neq x^*$, the *equilibrium condition*

$$(i) \qquad \pi(x^*, x^*) > \pi(x, x^*)$$

or, in the case of equity, the *stability condition*

$$(ii) \qquad \pi(x^*, x^*) = \pi(x, x^*), \pi(x, x) < \pi(x^*, x)$$

holds.

Note that $\pi(x, \widehat{x})$ is the payoff of the first player with strategy $x$ in a game with a second player playing strategy $\widehat{x}$.

The basic aim of ESS is to conceptualize requirements that capture the idea of a *stable* population-strategy by avoiding calculations of the complicated dynamics that are naturally entailed by evolutionary selection processes. Hence, the evolutionary quality of ESS is *static* in the sense that ESS does not explain *how* a population reaches an evolutionarily stable equilibrium state, but regards the question whether a population, having reached such a state, is apt to prevent alternative strategies, in sufficiently low frequency, from invading.

The fundamental *dynamical* conception of evolutionary game theory outlines evolutionary selection via the replicator dynamics (see Taylor and Jonker, 1978), determined by the system of ordinary differential equations (also known as the *replicator equation*),

$$\dot{p}_i(t) = [\pi(e_i, P(t)) - \pi(P(t), P(t))]p_i(t), \qquad (2.1)$$

with a dot "·" symbolizing the derivative with respect to time $t$. For finite-strategy games, the population state at moment $t$ is given by the finite-dimensional vector $P(t) = (p_1(t), ..., p_n(t))$ where $p_i(t)$ is the proportion of the players using pure strategy $e_i$ at that instant, such that $\sum p_i(t) = 1$ and $p_i(t) \geq 0$ hold true. The expected payoff of an individual with strategy-type $i$ is $\pi(e_i, P(t)) = \sum_{j=1}^{n} \pi(e_i, e_j) p_j(t)$

and the average payoff of the population (known as the *population fitness* which is mathematically equal to the expected payoff of a randomly chosen individual) is $\pi(P(t), P(t)) = \sum_{j=1}^{n} \pi(e_j, P(t)) p_j(t)$. Thus, from Eq. (2.1), the principle of the replicator dynamics is to benefit strategies that are above *population fitness* by spreading with a rate that is proportional to resulting fitness and to diminish those with a lower than average fitness in the same way.

There is an important link between ESS and dynamic stability of the finite replicator dynamics which brings the two basic properties of evolutionary selection processes, namely mutation (ESS) and selection (the replicator dynamics), together. The link is sufficiently clarified in an unambiguous manner: ESS is sufficient (though not necessary) for *asymptotic stability*[8] with respect to the replicator dynamics in population games with finite strategy set $S = \{1, ..., n\}$, also known as *matrix-games*, in which the distance of two populations is given in the natural Euclidean space (cf. Hofbauer et al., 1979).

### 2.2.2 Continuous Strategy Space and the Issue of Closeness

However, proving dynamic stability with respect to continuous strategy spaces is more challenging. Following the standard formulations in this regard, a population is described by a probability measure $P$ and identified with the aggregate play of its members on the measure space $(S, \mathscr{B})$ where $\mathscr{B}$ denotes the Borel $\sigma$-algebra on the compact metric space of feasible strategies, $S = [\underline{s}, \overline{s}]$.

$\pi : S \times S \rightarrow \mathbb{R}$ is assumed to be a bounded and Borel measureable payoff (or fitness) function. By virtue of compactness of $S$, the simplex $\Delta(S)$ is compact and metrizable (Oechssler and Riedel, 2002), and locates all populations on $S$. In this setting, a slightly different notation of ESS, for *Evolutionarily Stable State*, is more

---

[8]Intuitively, a population state $P$ is *asymptotically stable* if any sufficiently small change of the population composition results in a back drift toward $P$. A formal description of *asymptotic stability* (and the weaker criterion of *Lyapunov stability*) is given with Definition 2.2 below.

functional.

**Definition 2.1.** *A population $P$ is called an evolutionarily stable state (ESS) if for "mutation" $Q$, there is an invasion barrier $\epsilon(Q) > 0$ such that for all $0 < \eta \leq \epsilon$*

$$E(P, (1 - \eta)P + \eta Q) > E(Q, (1 - \eta)P + \eta Q). \tag{2.2}$$

Note that $E(P, \widehat{P})$ is the average payoff of population $P$ against the rival population $\widehat{P}$.

Strategies are now typically confined to pureness such that a personal mixed strategy "converts" to a *heterogeneous* population state, i.e. at least two different strategies are present in the population.[9] In this setting, ESSs are no longer necessarily asymptotically stable in the continuous replicator dynamics, given by

$$\dot{P}(t)(A) = \int_A \sigma(x, P(t))P(t)(dx), \tag{2.3}$$

with arbitrary subset $A \in \mathscr{B}$ of $S$. If $A$ describes a single point on the real line, then we have the standard replicator dynamics with a finite strategy space. The *differential fitness* of pure strategy $x$ playing against population $P$ is now given by

$$
\begin{aligned}
\sigma(x, P(t)) : \quad &= \quad E(\delta_x, P(t)) - E(P(t), P(t)) \\
&= \quad \int_S \pi(x, y)P(t)(dy) - \int_S \int_S \pi(x, y)P(t)(dy)P(t)(dx),
\end{aligned}
$$

with $\delta_x$ denoting the Dirac delta distribution of the homogeneous population with unit mass on $\{x\}$ (i.e. all present individuals exhibit the same strategy $x$).

Furthermore, even the stronger requirements of a strict Nash equilibrium, in which the equilibrium condition of ESS holds, are not sufficient to guarantee the likewise weaker requirements of Lyapunov stability. Whereas even a neutrally stable state (NSS), in which the ESS condition (Ineq. (2.2)) is allowed for equity, is sufficient for Lyapunov stability in the finite case (cf. Weibull, 1995, section 3.5).

---

[9]The original replicator dynamics of Taylor and Jonker (1978) treats personal strategies as *pure* determinants, but see for example Bomze (1991) for replicator dynamics with personal *mixed* strategies.

Oechssler and Riedel (2002) show this key result by example (see their Example 1 with fitness function $f(x, y) = -x^2 + 4xy$, where the Dirac delta $\delta_0$ describes a strict Nash equilibrium, and thus an ESS, but is not Lyapunov stable).[10]

The following definition specifies dynamic stability of the replicator dynamics.

**Definition 2.2.** *Let $Q^*$ be a rest point of the replicator dynamics.[11] Then*

- *$Q^*$ is called Lyapunov stable if for all $\epsilon > 0$ there exists an $\eta > 0$ such that*

  $\|Q(0) - Q^*\| < \eta \Rightarrow \|Q(t) - Q^*\| < \epsilon$ *for all $t > 0$.*

- *$Q^*$ is called asymptotically stable if additionally there exists $\epsilon > 0$ such that*

  $\|Q(0) - Q^*\| < \epsilon \Rightarrow \|Q(t) - Q^*\| \to 0$.

However, Bomze (1990) proves that the requirements of *strong uninvadability* are sufficient to guarantee dynamic stability of the replicator dynamics if the variational norm is used. Population $P$ is *strongly uninvadable* if there is an invasion barrier $\epsilon > 0$ such that $E(P, Q) > E(Q, Q)$ for all populations $P \neq Q$ with distance $0 < \|P - Q\| \leq \epsilon$. Oechssler and Riedel (2001) extend this finding by assuming a homogeneous population $P$ which is *uninvadable*, namely that conform to the ESS requirements of Definition 2.1 but with a uniform invasion barrier $\epsilon > 0$. This result is due to the fact that the criteria of uninvadabilty and strong uninvadability coincide in the case of homogeneous population play.

Recall, however, that in these articles the adopted topology is determined by the variational norm where the distance between two populations $P$ and $Q$ is given by

$$\|P - Q\| = 2\sup_{A \in \mathscr{B}} |P(A) - Q(A)|,$$

---

[10]Hofbauer et al. (2009) show that this result holds likewise for the BNN dynamics where new strategies can emerge if they yield better than average payoff. This happening is contrary to the replicator dynamics where the support is invariant at all times (see the next footnote for a definition of "support").

[11]A rest point $Q^*$ of the replicator dynamics is a Nash equilibrium with exactly the support of $Q^*$ as pure strategies. The support of a population $P \in \Delta(S)$ is the unique (relatively) closed subset of $S$ whose complement has measure 0 (with respect to $P$) and every open set that intersects it has positive measure.

(see Shiryaev, 1995, p. 360).

Yet, there is one major difficulty with respect to the informative values of studies using this definition of closeness: If a large fraction of a population changes from the equilibrium strategy then the new population is no longer close to the previous one even if the change is very small on the real line, thus potentially insignificant for the analysis at hand. For example, a population $P$ is $\epsilon$-close to another population with Dirac delta $\delta_y$ only if $\|P - \delta_y\| = \int |P - \delta_y| \, dx = 2\left(1 - P\left(\{y\}\right)\right) \leq \epsilon$. Hence, the fraction of pure $y$-players in population $P$ must conform to $P\left(\{y\}\right) \geq 1 - \frac{\epsilon}{2}$, which is a very strong barrier. To carry this example further, the distance of a homogeneous population is always maximal to a different homogeneous population even if the corresponding strategies are highly similar.

Oechssler and Riedel (2002) argue that the weak topology is generally the more appropriate alternative to measure the distance between populations, primarily because the weak topology regards the natural resemblance of slightly different strategies.[12] Formally, $P\left(t\right)$ converges weakly to $P^*$ if and only if $\lim_{t \to \infty} \int_S \pi dP\left(t\right) = \int_S \pi dP^*$ for every bounded, continuous real function $\pi$. If the Prohorov metric is used, then the distance of two populations is given by (cf. Billingsley, 1968, p. 238),

$$\rho\left(P, Q\right) = \inf\left\{\epsilon > 0, Q\left(A\right) \leq P\left(A^\epsilon\right) + \epsilon \text{ and } P\left(A\right) \leq Q\left(A^\epsilon\right) + \epsilon, \forall A \in \mathscr{B}\right\}$$

where $A^\epsilon = \left\{x : \exists y \in A, |y - x| < \epsilon\right\}$.

Thus, for the sake of closeness, it is irrelevant whether a population changes such that a small fraction plays a very different strategy or if a large fraction changes such that a very likely strategy is played. To make this conclusive, think of the following situation. If there is a two-type population $P$ that consists of $y$-players and $\epsilon \in [0, 1]$-fractional of $x$-players, then the distance between population $P$ and the

---

[12]Apparently, this is generally a proper reason in favor of the weak topology but there are also some arguments which promote the strong topology in some settings. One can think of situations where an initial population should not be close to the mutated new population even if the bulk of the former changes only very little on the real line since this could also be a significant evolutionary incident (*every* individuum of the bulk changes his or her strategy), which should be reflected by the topology at hand. Another, more technical, reason in favor of the strong topology is that this topology does not require continuity of the underlying fitness function.

Dirac measure on $\{x\}$ is $\rho(P, \delta_x) = \min\{\epsilon, |y - x|\}$. Moreover, two homogeneous populations are close to each other if the corresponding pure strategies are close on the real line in the Euclidean topology. By virtue of these arguments, we think that the weak topology measures best the nature of evolution of preferences.

**Remark.** *Throughout this paper we assume closeness and convergence with respect to the weak topology (unless otherwise stated).*

### 2.2.3 CSS, NIS, and $\mathscr{ER}$

Intuitively, ESS is insufficient for dynamic stability in the continuous strategy case since this concept is unique to the situation that a *small* fraction of the population changes from an established strategy. From this perspective, the biological concepts CSS (see Eshel and Motro, 1981; Eshel, 1983) and NIS (see Apaloo, 1997) enlarge the stability analysis by regarding resistance against a mutated bulk. Both concepts are originally analyzed for an ecological system in which the feasible strategy set evolves according to a one-dimensional continuity of pure alternatives. This appearence is in line with our detection of the evolution of altruism and envy assumed to be opposed preferences in a one-dimensional continuity. Oechssler and Riedel (2002) have introduced the strong criterion $\mathscr{ER}$ to unify the characteristics of these concepts and to provide a strong argumentation for evolutionary fitness in a wide variety. By regarding the essence of the weak topology, $\mathscr{ER}$ is given if both kinds of mutations are unsuccessful, "a *large* change of strategic play by a *small* fraction of players as well as a *small* change of strategic play by a *large* fraction of the population".[13]

In what follows, we will firstly outline the formal definitions of CSS, NIS, and $\mathscr{ER}$, and then center on some consequences as discussed in the literature.[14]

---

[13]In Cressman and Hofbauer (2005), $\mathscr{ER}$ is called *locally superior* (with respect to the weak topology).

[14]Eshel and Sansone (2003) have introduced "Continuous Replicator Stability" (which is much alike to the conditions of NIS) to assess dynamic stability of the replicator dynamics. Though, this approach is only sufficient if the maximal shift topology is used which is a limitation of the

To further characterize CSS, assume in the continuous setting that the entire (homogeneous) population $\delta_{x^*}$ plays an ESS and then switches slightly from it by playing a very similar strategy $x = x^* + \omega$ with $0 < |\omega| < \epsilon$. This image gives rise to an important differentiation of two unlike types of ESSs with respect to their capability of surviving evolutionary pressures: ESSs that are *continuously stable*, and those that are not. Any ESS is CSS if a sufficiently small change of a population as a whole is such that mutations closer to the ESS are able to invade the new population via higher fitnesses, hence leading the population back in the direction of the ESS (though not necessarily to the ESS itself). To define CSS rigorously, consider the following.

**Definition 2.3.** *Any strategy $x^*$ is CSS if (1) it is ESS and (2) there exists an $\epsilon > 0$ such that for all $x$ with $|x^* - x| < \epsilon$ there exists a $\delta > 0$ such that for all $y$ with $|x - y| < \delta$*

$$\pi(y, x) > \pi(x, x) \quad \text{if and only if} \quad |y - x^*| < |x - x^*| .$$

Accordingly, the criterion of CSS can be divided into two matters. The first one is static in that any CSS is ESS. The second one gives an intuitive dynamic justification since the higher fitnesses of equilibrium-closer strategies drive the population toward the equilibrium. This extra condition derives from the adaptive dynamics approach where it is termed *m-stability* by Taylor (1989) and *convergence stability* by Christiansen (1991). For the adaptive dynamics approach, an interior CSS is an asymptotically stable rest point (e.g. Cressman, 2009).[15]

Now, we focus on more operable conditions for CSS which we continue to use for the game analysis in section 2.4. Consider the first and second order derivatives of

---

weak topology in that it requires additionally that the support of two populations have to be close to guarantee closeness (see also the remark in section 5 of Cressman et al., 2006). Thereof, we prescind from this approach.

[15]The population mean $\overline{x}$ evolves according to the *canonical equation* of adaptive dynamics, $\dot{\overline{x}} = g(\overline{x}) \pi_x(\overline{x}, \overline{x})$ where $g(\overline{x})$ is a positive function that is related to the rate mutations occur and the variance of their distribution.

the payoff function $\pi\left(x, y\right)$, namely $\pi_x$, $\pi_{xx}$, $\pi_{xy}$, at $x = y = x^*$ and the following conditions.

$$\pi_x\left(x^*, x^*\right) = 0 \tag{2.4}$$

and

$$\pi_{xx}\left(x^*, x^*\right) \leq 0. \tag{2.5}$$

With Eq. (2.4) the strict case of Ineq. (2.5) is sufficient for ESS. As pointed out by Eshel (1983), a necessary condition for any ESS $x^*$ to be CSS is given whenever

$$\pi_{xx}\left(x^*, x^*\right) + \pi_{xy}\left(x^*, x^*\right) \leq 0 \tag{2.6}$$

holds. The strict case of Ineq. (2.6) is sufficient for CSS.

From this specification, Eshel (1983) proposes a geometrical interpretation of CSS. Any optimal strategy $x$ played against strategy $y$ can be expressed by a *best reply function* $x\left(y\right) = x$ such that $\pi_x\left(x\left(y\right), y\right) = 0$ is necessary. By implicit differentiation, the first order derivative is given by

$$\pi_{xx}\left(x\left(y\right), y\right) \cdot x'\left(y\right) + \pi_{xy}\left(x\left(y\right), y\right) = 0.$$

Rearranging gives

$$x'\left(y\right) = -\frac{\pi_{xy}\left(x\left(y\right), y\right)}{\pi_{xx}\left(x\left(y\right), y\right)}. \tag{2.7}$$

Eq. (2.7) together with the sufficient conditions of ESS and CSS state

$$\pi_{xx}\left(x^*, x^*\right) + \pi_{xy}\left(x^*, x^*\right) < 0 \;\Leftrightarrow\; 1 + \frac{\pi_{xy}\left(x\left(y\right), y\right)}{\pi_{xx}\left(x\left(y\right), y\right)} > 0,$$

and hence, CSS is guaranteed if

$$x'\left(y\right) = -\frac{\pi_{xy}\left(x\left(y\right), y\right)}{\pi_{xx}\left(x\left(y\right), y\right)} < 1. \tag{2.8}$$

Any intersection of the best reply function and the main diagonal $x = y$ is obviously ESS since any strict best reply to itself is ESS, which is pursuant to the equilibrium condition. According to Ineq. (2.8), any intersection is CSS, if the best reply function intersects additionally 'from above'.

A NIS (also named a *good invader* by Kisdi and Meszéna, 1995) exposes the features of a strategy that is able to invade any strategy that is sufficiently close (with respect to the Euclidean norm). However, NIS does not require ESS but a NIS that holds for ESS is called ESNIS for *Evolutionarily Stable Neighborhood Invader Strategy* (see Apaloo, 2005). A CSS is eventually not apt to invade a close neighbor. Hence, through evolutionary dynamics it may be repelled by a very similar strategy which is an incident that cannot happen to an ESNIS. Any ESNIS is a CSS but the converse cannot be taken for granted in general.

A formal definition of NIS is given with the following.

**Definition 2.4.** *Any strategy $x^*$ is NIS if $x^*$ can invade any $x \neq x^*$ in an $\epsilon > 0$ neighborhood, i.e.*

$$\pi(x^*, x) > \pi(x, x), \forall x \quad with \ |x - x^*| < \epsilon.$$

Oechssler and Riedel (2002) show that a necessary condition is given by

$$\pi_{xx}(x^*, x^*) + 2\pi_{xy}(x^*, x^*) \leq 0. \tag{2.9}$$

Analogue to the CSS condition, the strict case of Ineq. (2.9) is sufficient for NIS. The stronger conditions of ESNIS state geometrical interpretations that are similar to the CSS ones. We summarize these useful facts in the following Lemma.

**Lemma 2.1.** 1) Any ESS $x^*$ is a CSS, if the best reply function $x(y)$ intersects the main diagonal $x = y$ at $x^*$ from above. 2) Any ESS $x^*$ is an ESNIS, if the best reply function $x(y)$ intersects the main diagonal $x = y$ at $x^*$ with a slope smaller than $\frac{1}{2}$.

*Proof.* The first statement is due to Eshel (1983, Theorem 2). However, the proof can also be retraced with the explanations above. The proof of the second statement is similar. By regarding the sufficient condition of ESS, the sufficient condition of NIS states

$$-\frac{\pi_{xy}(x(y), y)}{\pi_{xx}(x(y), y)} < \frac{1}{2}.$$

26

We know that $x'(y) = -\frac{\pi_{xy}(x(y),y)}{\pi_{xx}(x(y),y)}$, such that $x'(y) < \frac{1}{2}$ is deciding. $\qquad\square$

The strong concept of Oechssler and Riedel (2002) is defined as follows.

**Definition 2.5.** *A population $P^*$ is evolutionary robust if there exists $\epsilon > 0$ such that for all $Q \neq P^*$ with $\rho(P^*, Q) < \epsilon$ we have $E(P^*, Q) > E(Q, Q)$.*

$\mathscr{ER}$ corresponds to *strong uninvadability* in the variational norm. Though, $\mathscr{ER}$ is much harder to attain, which is due to the larger set of mutations $Q$ that are $\epsilon$-close to $P^*$ in the weak topology. An equilibrium which is $\mathscr{ER}$ guarantees dynamic stability for the replicator dynamics in doubly symmetric games (games in which the players' payoff always coincide). However, Oechssler and Riedel (2002) only conjecture that $\mathscr{ER}$ is sufficient for dynamic stability in the weak topology for general $\pi(x, y)$. Furthermore, there is no simple solution algorithm for checking $\mathscr{ER}$-strategies, but an approach with quadratic payoff functions is given in the Appendix.

Further implications of CSS, NIS, and $\mathscr{ER}$ in (evolutionary) game-theoretical meanings are widely discussed (see the recent papers of Eshel and Sansone, 2003; Apaloo, 2005; Cressman, 2005; Cressman and Hofbauer, 2005; Cressman et al., 2006; Hofbauer et al., 2009; Cressman, 2009). For example, Cressman (2009) relates CSS and NIS to classical game-theoretic solution concepts when applied to two-player games with a continuous strategy space. A CSS $x^*$ in the interior of the strategy space is equivalent to *neighborhood $\frac{1}{2}$-superiority* (i.e., $\pi(x^*, P) > \pi(P, P)$ for all $P \in \Delta(S)$ with $1 > P(\{x^*\}) \geq \frac{1}{2}$ with support sufficiently close to $x^*$).[16] Moreover, a neighborhood strict Nash equilibrium (cf. the *equilibrium condition* of ESS with $x^*$ close to $x$) which is NIS is equivalent to *neighborhood $0$-superiority*. Also, Cressman (2009) identifies a dynamical consequence of CSS by considering the Cournot adjustment process of Moulin (1984) with an interior CSS as an asymptotically stable rest point. Hofbauer et al. (2009) show the relevance of CSS and $\mathscr{ER}$ in several

---

[16]Strategy $x^*$ is *globally $\frac{1}{2}$-superior* if $\pi(x^*, P) > \pi(P, P)$ for all $P \in \Delta(S)$ with $1 > P(\{x^*\}) \geq \frac{1}{2}$.

dynamics. For quadratic fitness functions of the form $\pi(x,y) = -x^2 + axy$, they give a thorough analysis of the replicator dynamics, the BNN dynamics, and the best response dynamics (see Table 1 in chapter 8 there). Most notably, a CSS is asymptotically stable for the BNN dynamics and the best response dynamics. A strategy which achieves the stronger conditions of $\mathscr{ER}$ is asymptotically stable for the replicator dynamics. Cressman (2005) shows some consequences of CSS and NIS for dynamic stability of the replicator dynamics which are unambiguous and extensive if $\pi(x,y) = \pi(y,x)$. However, the connections with respect to arbitrary payoffs are much more tenuous. In this regard, Eshel and Sansone (2003) give the NIS concept a strong meaning since they show that NIS is a necessary condition for dynamic stability.

Some further useful implications of CSS and NIS are due to Cressman and Hofbauer (2005) and their measure dynamics approach in the one-dimensional continuum. They show that a homogeneous population which is not CSS cannot be dynamically stable, but a CSS is dynamically stable if the initial distribution is close to the CSS and the support is a compact interval with the CSS in the interior. A somewhat stronger meaning holds for a homogeneous population which is NIS since an arbitrary initial support (however still close to the equilibrium) is then already sufficient.

## 2.3   Indirect Evolution

In this section we describe the approach of indirect evolution in detail. By doing so, we introduce some game notations and assumptions that are valid for the remainder of this paper.

Let $\Gamma(S, T, \pi_1(x,y), \pi_2(x,y), U_1(x,y), U_2(x,y))$ denote the evolutionary symmetric 2-person game in the conventional setting of one infinitely large population in which individuals are drawn randomly (with equal probability) and repeatedly to

play a certain game at each stage.[17] Sticking to the rationality-acceptation of ortho-dox game theory, the two players, identified with subscripts on the game-functions, are in the position to act perfectly according to their subjective preferences by max-imizing their well-being $U_1(x, y)$ and $U_2(x, y)$ with equilibrium strategies $x^* \in S$ (player 1) and $y^* \in S$ (player 2). Contrary to models of standard *direct* evolu-tion, the functions that players seek to maximize are no longer equivalent to their monetary payoffs, $\pi_1(x, y)$ and $\pi_2(x, y)$, which, however, drive their evolutionary success.[18] Fixing evolutionary fitness with equilibrium behavior leads to the new fitness function with evolving preference-types $\alpha \in T$ (player 1) and $\beta \in T$ (player 2). The fitness of player 1 is then given by $f_1(\alpha, \beta) = \pi_1(x^*(\alpha, \beta), y^*(\alpha, \beta))$.[19] Hav-ing these results allows the researcher to continue in analogy to standard direct evolution but with evolving intrinsic, *subjective* values on trial. For example, Nash equilibria *preferences* can be found at equilibrium behavior with the usual definition $f(\alpha, \beta) \leq f(\alpha, \alpha)$ for all different $\beta$. This technique holds likewise for all refinement concepts on study.

Consequently, preferences that are sufficiently fit are more reproductive than less fit preferences, such that a final population state consists of preferences gaining most, though not necessarily maximizing (directly), evolutionary reproductive suc-cess.

Notice that indirect fitness-detection of preferences combines traditional assump-tions of *rational* decision making, like in orthodox game theory, with evolutionary adapting concepts resulting from evolutionary game theory.

---

[17]By the *infinity* assumption, we can exclude further random effects which would result in stochastic evolutionary dynamics as studied, e.g., in Güth et al. (2002).

[18]By virtue of the economic context we imply to the indirect evolutionary game and for the sake of convenience, we will assume that monetary profit is the basic determinant for reproductive success. In fact, a couple of studies show that economic success is closely related to the number of surviving offspring (e.g. Boyer, 1989) which legitimate the replicator dynamics as the selection process. But note that it is only important that the fitness criterion is an *objective* (or rather intersubjective ackknowledged) value which represents success in the society. For example, a higher status (without earning more money) may indicate a similar relevance.

[19]Since the position of the players is by no means the equilibrium fitness of player 2 can easily be found by symmetry. This holds likewise for later analysis.

The basic steps of indirect evolution of preferences are illustrated in Figure 2.1.



Figure 2.1: One Sequence of Indirect Evolution of Preferences

*To avoid fundamental questions about the 'genesis' of individual preferences, we restrict our attention to an evolutionary process which is already in progress. At step ①, individuals play rationally (and immediately) for their perceptive well-being which possibly subsumes an importance on others payoff. The equilibrium play then fix the objective outcome (which is the evolutionary success) at the ②nd step. At step ③, the objective outcome determines the evolutionary selection since the underlying evolutionary process favors a higher outcome by spreading. Preferences that are sufficiently fit are bequeathed from one generation to another in an accordant high continuance and generate the successive population in evolutionary time, which is the ④th step.*

The indirect evolutionary approach has become the central reply to (experimental) evidence that dispute to homo economicus models with the basic postulate of material selfishness. However, indirect preferences evolution cannot replace exper-

imental studies of preferences sometimes changing in very short-termed experimental situations since evolution of preferences is a long-termed process appraising how individuals *tend* to act in certain strategic situations.

Since we limit ourselves to the replicator dynamics as the most right choice to model the preference shape, we restrict the validity of indirect evolution to selection in a genetical sense. But indirect evolution continues its importance, e.g., in social learning processes where preferences are shaped by adaption from "role models" (like parents or teachers).

## 2.4   The Model with CSS, NIS, and $\mathscr{ER}$

The basic formulation of our evolutionary game is similar to Bester and Güth (1998), but incorporates the suggestions of Bolle (2000), and, more explicitly, of Possajennikov (2000).

Analogous to the notations above, let

$$\pi_1(x, y) = x(ky + m - x) \text{ and } \pi_2(x, y) = y(kx + m - y) \tag{2.10}$$

denote the players' monetary payoffs, fixed by the strategy-choices $x, y \in S = \mathbb{R}_0^+$ and the parameters $m > 0$ and $k$ which is either positive or negative. These payoffs can represent several strategic situations. For example, one can think of a simple production game with efforts $x$ and $y$ where the players either exploit a resource or contribute to a public good. The specific situation then depends on the type of externality. Alternatively, the payoffs may represent profit functions in a symmetrical duopoly competition with heterogeneous products and linear demand. Accordingly, in a Bertrand market the choices $x$ and $y$ would define the prices while in a Cournot market $x$ and $y$ would be the quantity choices. The parameter $k$ (and especially the *sign* of $k$) represents a basic characterization of the game. At first we follow Bester and Güth that $-1 < k < 1$; however, the $k$ is further specified

31

below. Note that $k$ determines interdependency in the following manner: A positive $k$ leads to $\frac{\partial \pi_1(x,y)}{\partial y} > 0$ which means that a higher input of the respective co-player implies a higher fitness-level, i.e. the players impose positive externalities on one another. The opposite holds for $k < 0$, since in this case $\frac{\partial \pi_1(x,y)}{\partial y} < 0$ and negative externalities are present.

The players derive evolutionary success from their monetary profit, but they are foremost provided to maximize their subjective well-being, given by

$$U_1(x,y) = \pi_1(x,y) + \alpha \pi_2(x,y) \text{ and } U_2(x,y) = \pi_2(x,y) + \beta \pi_1(x,y). \qquad (2.11)$$

Unless a co-player's profit is irrelevant, meaning a player's attitude is egoism ($\alpha$ or $\beta = 0$), the players have altruistic ($\alpha$ or $\beta > 0$) or envious ($\alpha$ or $\beta < 0$) preferences on others profit. To avoid preferences that go somewhat beyond envy or altruism such that the players care less about themselves than about others, we restrict the type-parameters to the region $\alpha, \beta \in T = [-1, 1]$.[20]

Recall that the players are equipped with the capability of perfect observation of co-players preferences. Hence, the context of our game is in situations where individuals either know each other very well or learn preferences sufficiently fast.[21]

Following the indirect evolutionary approach, the players are in the position to maximize their well-being. Formally

$$x^* \in \underset{x}{\operatorname{argmax}}\, U_1(x, y^*), \quad y^* \in \underset{y}{\operatorname{argmax}}\, U_2(x^*, y), \qquad (2.12)$$

with resulting *reaction functions*

$$x = \frac{k(\alpha + 1)y + m}{2}, \quad y = \frac{k(\beta + 1)x + m}{2}. \qquad (2.13)$$

According to the terminology of Bulow et al. (1985), strategies are *complements* to each other with $k > 0$ since the slope of the reaction function is then positive,

---

[20]Possajennikov and Bolle relax assumptions on game-parameters by allowing for wider ranges. But this does not affect qualitative results of the region stated here.

[21]Retrace the argumentation of Frank (1987) for physical indications help explaining why this is not only an artificial assumption.

meaning that a higher input of the co-player leads to a higher marginal revenue. And strategic *substitutes* are determined by $k < 0$ with consistencies the other way round.

The intersection point of the reaction functions represents equilibrium-play[22]

$$x^*(\alpha, \ \beta) = \frac{m(k(\alpha + 1) + 2)}{4 - k^2(\alpha + 1)(\beta + 1)}, \quad y^*(\alpha, \ \beta) = \frac{m(k(\beta + 1) + 2)}{4 - k^2(\alpha + 1)(\beta + 1)} \quad (2.14)$$

which makes the impact of player A's preference $\alpha$ (respective $\beta$) on player B's behavior $y^*$ (respective $x^*$) transparent.

By fixing the monetary payoff with the equilibrium strategies, one reaches the new fitness function of player 1 with evolving preferences:

$$
\begin{aligned}
\pi_1(x^*(\alpha, \beta), y^*(\alpha, \beta)) &= f_1(\alpha, \ \beta) \\
&= -\frac{m^2(k(\alpha + 1) + 2)(k^2\alpha(\beta + 1) + k(\alpha - 1) - 2)}{(4 - k^2(\alpha + 1)(\beta + 1))^2} \quad (2.15)
\end{aligned}
$$

while the fitness function of player 2 satisfies $f_2(\beta, \ \alpha) = f_1(\alpha, \ \beta)$. Due to the indirect evolutionary approach, the ESS conditions of section 2.2.1 can now be transferred to preferences such that

(i)  $f_1(\alpha^*, \ \alpha^*) > f_1(\alpha, \ \alpha^*) \ \forall \ \alpha \in [-1, 1]$, or

(ii)  $f_1(\alpha^*, \ \alpha^*) = f_1(\alpha, \ \alpha^*) \quad$ and $\quad f_1(\alpha^*, \ \alpha) > f_1(\alpha^*, \ \alpha) \ \forall \ \alpha \in [-1, 1]$

is deciding for $\alpha^*$ to be ESS.

Calculating the first order condition yields

$$\alpha = -\frac{k(\beta + 1)(k + 2)}{\beta k(k - 2) + k^2 - 2k - 4}. \quad (2.16)$$

Any preference $\alpha^*$ is an ESS candidate if it is a best reply to itself. Hence, setting $\alpha = \beta = \alpha^*$ and solving for the equilibrium preference guides to ESS-candidates: $\alpha_1^* = -\frac{k+2}{2}$, $\alpha_2^* = \frac{k}{2-k}$, and corner-solutions $\alpha_3^* = -1$, and $\alpha_4^* = 1$. Straightforward calculations state the following proposition which is due to Possajennikov (2000) for a wider range of $k$, namely $-2 \leq k < 1$, $k > 2$, and existing strategic interdependence ($k \neq 0$), and to a former analysis of Possajennikov (1999) with the same

---

[22]Eqs. (2.14-2.16) coincide with Eqs. (3-5) in Possajennikov (2000).

constraints on $k$.

**Proposition 2.1.** $\alpha^* = \frac{k}{2-k}$ is the unique ESS preference.

Since $k > 0$ results in $\alpha^* > 0$, it is verified that altruism is ESS with strategic complements. On the other hand, strategic substitutes leads to envy being ESS because $\alpha^* < 0$ if $k < 0$.[23] Note that with the restriction of $k$, ESS preferences are always in the interior of the feasible space. An interesting aspect arises by exploring the effect of the strategic situation on the strength of the evolutionarily stable preference. If strategic complements approaches 1, then the evolutionary players are nearly perfectly altruistic in equilibrium, while the situation of strategic substitutes determines at most one-third of the whole range of envy. Figure 2.2 illustrates the relation of strategical interdependence on ESS preferences.

**Remark.** *Heifetz et al. (2007a) have analyzed the same setting with some further restrictions on k (or "−b" in their case; although they have used somewhat different notations on the parameters, the game is identical), namely $-\frac{2}{5} < k < \frac{1}{2}$ and $k \neq 0$, to explore whether evolutionary dynamics will evolve to the ESS preference. They have shown that any initial population with full support on the preferences space $\alpha, \beta \in [-1, 1]$ will evolve to the evolutionarily stable state under any payoffmonotonic and regular selection dynamics.[24] Recall that this preference space is the same that we assume. As a limitation of k, they have proven that Eq. (2.15) is strictly concave in the preferences only if $k > -\frac{1}{2}$. In what follows, we adopt this assumption such that we can work with a best reply function.*

Since the concept of ESS is too weak in the continuous space of preferences,

---

[23]This basic finding is restricted to the more natural cases of $-1 < k < 1$ where own action has got more impact on outcome than the action of the opponent. With respect to the wider parameter range of Possajennikov (2000), an $\alpha^* < 0$ with an absolute value larger than one can even be ESS if strategies are complements with sufficient high interdependence ($k > 2$).

[24]Dynamics are payoff monotone if a higher average fitness corresponds to a higher growth rate, or formally $\frac{1}{P(A)} \int_A E(\delta_x, P) P(dx) > \frac{1}{P(A')} \int_{A'} E(\delta_x, P) P(dx) \Leftrightarrow \frac{\dot{P}(A)}{P(A)} > \frac{\dot{P}(A')}{P(A')}$. Dynamics are regular if evolution does not allow for innovations which means that $\dot{P}(A) = 0$ holds for all $A \in \mathscr{B}$ with $P(A) = 0$.

Figure 2.2: Relation of strategical interdependence and ESS preferences.

we proceed by using the stability concepts of section 2.2.3. At first, we reach the following result.

**Proposition 2.2.** The strict Nash equilibrium preference $\alpha^* = \frac{k}{2-k}$ is CSS.

*Proof.* From section 2.2.3 we know that any intersection of the main diagonal $x = y$ with the best reply function $x(y)$ is an ESS. Any ESS $x^*$ is additionally a CSS if the best reply function intersects the main diagonal *from above.* Following the indirect evolutionary approach, this condition can be transferred to our preference setting. Eq. (2.16) is now considered as of the first players' best reply function $h(\beta)$,

$$h(\beta) = \alpha = -\frac{k(\beta+1)(k+2)}{\beta k(k-2) + k^2 - 2k - 4}. \tag{2.17}$$

It is left to check whether the best reply function intersects the main diagonal at $\alpha^* = \frac{k}{2-k}$ with a slope smaller than 1. The slope of the best reply function at the equilibrium preference $\alpha^* = \alpha = \beta = \frac{k}{2-k}$ is

$$\left.\frac{dh(\beta)}{d\beta}\right\}_{\beta=\frac{k}{2-k}} = \frac{k}{k+2}. \tag{2.18}$$

Eq. (2.18) always underbids the declared value of 1 since with the restrictions of $k$

we have

$$\frac{dh(\beta)}{d\beta}\bigg\}_{\beta=\frac{k}{2-k}} = \frac{k}{k+2} < 1 \quad \Rightarrow \quad 0 < 2. \tag{2.19}$$

Figure 2.3 shows that the best reply function intersects the main diagonal from above.

□

Now, we apply the NIS concept.

**Proposition 2.3.** The strict Nash equilibrium preference $\alpha^* = \frac{k}{2-k}$ is NIS.

*Proof.* A sufficient condition for NIS is given by Ineq. (2.9). Applying this condition to the relevant fitness function Eq. (2.15), the sufficient condition becomes

$$f_{\alpha\alpha}\left(\alpha^* = \frac{k}{2-k}, \alpha^* = \frac{k}{2-k}\right) + 2f_{\alpha\beta}\left(\alpha^* = \frac{k}{2-k}, \alpha^* = \frac{k}{2-k}\right) < 0. \tag{2.20}$$

After calculating, we get

$$f_{\alpha\alpha}\left(\alpha^* = \frac{k}{2-k}, \alpha^* = \frac{k}{2-k}\right) + 2f_{\alpha\beta}\left(\alpha^* = \frac{k}{2-k}, \alpha^* = \frac{k}{2-k}\right) = \frac{m^2k^2(k-2)^6}{512(k-1)^3}. \tag{2.21}$$

With the restrictions of $m$ and $k$ the result is always negative.

Furthermore, from Lemma 2.1 and the restrictions on $k$ it follows that

$$\frac{dh(\beta)}{d\beta}\bigg\}_{\beta=\frac{k}{2-k}} = \frac{k}{k+2} < \frac{1}{2} \quad \Rightarrow \quad k < 2 \tag{2.22}$$

is true, which is sufficient for NIS, too.

□

To show the following result, we translate evolutionary fitness to a quadratic function which allows us to use the classification scheme as proposed by Cressman and Hofbauer (2005).

Figure 2.3: Intersection of $h(\beta)$ with main diagonal $\alpha = \beta$.

**Proposition 2.4.** In the Taylor-approximated case of a quadratic fitness function, the strict Nash equilibrium preference $\alpha^* = \frac{k}{2-k}$ is $\mathscr{ER}$.

*Proof.* See the Appendix. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Hence, the homogeneous population with unit mass on $\left\{\alpha^* = \frac{k}{2-k}\right\}$ is both strong against small shifts by a large fraction of the population and strong against large shifts by a small fraction of the population. Recall, that $\mathscr{ER}$ implies CSS as well as NIS. Note also that $\mathscr{ER}$ implies (strong) uninvadability (see Oechssler and Riedel, 2002). The latter implication gives the replicator dynamics a strong meaning if the variational norm is used on the population space since a homogeneous, uninvadable population is then asymptotically stable with respect to the infinite preference space (see Oechssler and Riedel, 2001). However, as mentioned above, we basically assume the weak topology to explain the evolution of preferences. As discussed, identifying

37

the stability criteria is more subtle here, but $\mathscr{E}\mathscr{R}$ is a strong argument in favor of the equilibrium $\alpha^* = \frac{k}{2-k}$.

## 2.5 The Analysis with Sophisticated-Perceptive Preferences

So far, other-regarding preferences have been formulated with subjective well-being that is specified by a utility function predefined through one's own monetary payoff in addition to a certain extent of a co-player's monetary payoff. This section deals with the convention that individuals with other-regarding preferences take the *real* well-being of co-players into account. Thus, the perceived payoff function depends no longer on the monetary, objective payoff of others but on their subjective payoff which possibly differs from the objective payoff as well. This assumption is justified by the valid common knowledge assumption regarding preferences in addition to the standard premise of a *real* care in the relevant case.[25] Hence, the reformation of preferences implemented here describes a legitimate, more differentiated definition of altruism and envy, and individuals adherent with this definition exhibit a "sophisticated perception" of a co-player's well-being. This is why we suggest the term "sophisticated-perceptive preferences". The basic question is now then, whether this reformation makes a difference in the present evolutionary framework. To answer this question, let us first detail the new situation more formally.

We define the subjective well-being of the two players as follows.

$$V_1^n(x,y) = \pi_1(x,y) + \alpha V_2^{n-1}(x,y) \text{ and } V_2^n(x,y) = \pi_2(x,y) + \beta V_1^{n-1}(x,y), \quad (2.23)$$

where $n \in \mathbb{N}_0 = \{0,1,2,...\}$ is the number of perceptive iteration-steps which we assume is a homogeneous trait-value in the population. Besides, the special case

---

[25]Predominantly, this alternative definition is proper by virtue of the valid common knowledge assumption with respect to preferences but does not call for being more "true" in general since an accurate definition of altruism or envy seems always to depend heavily on the given environment and an underlying question one wants to answer.

of $n = -1$ requires an extra definition which follows a break down into monetary payoffs:

$$V_1^{n=-1}(x,y) = \pi_1(x,y) \text{ and } V_2^{n=-1}(x,y) = \pi_2(x,y).$$

All remaining game-parameters are the same as in section 2.4.

Now we are able to distinguish the different types of preferences with $n \in \mathbb{N}_0 \cup \{-1\}$. At first, the orthodox *homo economicus* assumption is valid if $n = -1$ since egoism is then represented. Other regarding preferences appear if $n \in \mathbb{N}_0$, and sophisticated-perceptive preferences are present if $n \in \mathbb{N} = \{1,2,3,...\}$. In particular, we are foremost interested in the case of $n \to \infty$, since we assume the ability to perfect perception, and the evolutionary differences to $n = 0$ (which is the case of the previous section with other-regarding preferences but without a sophisticated-perception).

With this issue in mind, we can now give the following result which shows the irrelevance of the existence of a sophisticated perception of a co-player's well-being.

**Proposition 2.5.** Consider the evolutionary environment discussed so far. Then, the existence of a sophisticated perception of well-being, which we assume is a homogeneous and "absolute" value,[26] is immediately negligible with respect to the evolutionary fitness of the equilibrium preference. Precisely, $x^* \in \underset{x}{\operatorname{argmax}} \lim_{n \to \infty} V_1^n(x,y^*) = x^* \in \underset{x}{\operatorname{argmax}} V_1^{n=0}(x,y^*) = x^* \in \underset{x}{\operatorname{argmax}} U_1(x,y^*)$.

*Proof.* We only give the proof with respect to the first players' well-being since the position of the players is inconsequential. By incorporating even and odd numbers $n$, the calculations of the limit values are given by

---

[26] By "absolute" value, we mean that the perceptive iteration-steps of the individuals run to infinity. Alternatively, we could consider subjective preference functions which are recursive in the sense of $U_1^\diamond(x,y) = \pi_1(x,y) + \alpha U_2^\diamond(x,y)$ (recall that the positions of the players is by no means such that for the second player $U_2^\diamond(x,y) = \pi_2(x,y) + \beta U_1^\diamond(x,y)$ holds) which would lead to the same qualitative finding of this proposition.

$$
\begin{aligned}
\lim_{n \to \infty} V_1^{2n+2} &= \lim_{n \to \infty} \left( 1 + \alpha\beta + .. + \alpha^{n+1}\beta^{n+1} \right) [\pi_1 + \alpha\pi_2] \\
&= \lim_{n \to \infty} \frac{1 - \alpha^{n+2}\beta^{n+2}}{1 - \alpha\beta} [\pi_1 + \alpha\pi_2] \\
&= \frac{1}{1 - \alpha\beta} [\pi_1 + \alpha\pi_2] \\
&= \frac{V_1^{n=0}}{1 - \alpha\beta},
\end{aligned}
$$

and

$$
\begin{aligned}
\lim_{n \to \infty} V_1^{2n+1} &= \lim_{n \to \infty} \left[ \left( 1 + \alpha\beta + .. + \alpha^{n+1}\beta^{n+1} \right) \pi_1 + \alpha \left( 1 + \alpha\beta + .. + \alpha^n\beta^n \right) \pi_2 \right] \\
&= \lim_{n \to \infty} \frac{1 - \alpha^{n+1}\beta^{n+1}}{1 - \alpha\beta} \left[ \pi_1 + \alpha^{n+1}\beta^{n+1}\pi_1 + \alpha\pi_2 \right] \\
&= \frac{1}{1 - \alpha\beta} [\pi_1 + \alpha\pi_2] \\
&= \frac{V_1^{n=0}}{1 - \alpha\beta}.
\end{aligned}
$$

Since there is only one limiting value for both cases, we can use this result.

A basic postulate of the indirect evolutionary approach that players seek to maximize their well-being implies that

$$
x^* \in \operatorname*{argmax}_x \frac{V_1^{n=0}(x, y^*)}{1 - \alpha\beta},
$$

which is equivalent to

$$
x^* \in \operatorname*{argmax}_x \frac{U_1(x, y^*)}{1 - \alpha\beta}.
$$

By a simple calculation, we see that

$$
x = \frac{k(\alpha + 1)y + m}{2}
$$

declares the first players' best reaction. Obviously, the denominator of the limit value is irrelevant for evolutionary fitness of the respective preference since the reaction function is the same as in the previous section (see Eq. 2.13), and the same computation under indirect evolution follows. $\qquad\square$

Although sophisticated-perceptive preferences appear like a somewhat artificial game-theoretical product, one can think of situations where human beings with other-regarding preferences anticipate the fact that others might feel similar. Consider, for example, the question of who will pay the bill in a restaurant with very close friends where one party has to give the other party's goodwill the precedence. However, we leave it to the reader to decide to what extent sophisticated-perceptive preferences play a relevant role in human society. Also, the question of what is a realistic number of perceptive iteration-steps $n$ is quite far beyond the scope of this study. Maybe this would be an interesting experiment to study in the laboratory—however, mounting such an experiment seems to be more difficult than the traditional ones with usual other-regarding preferences.

## 2.6 Conclusion

In this paper, we have reviewed the study of altruism by Bester and Güth (1998), and the one-dimensional enlargement of the preference space by Possajennikov (2000) and Bolle (2000), respectively. By making two expedient assumptions on the evolution of preferences, this study is a straightforward extension of the former surveys and furthermore a required task, as shown by Oechssler and Riedel (2002). If one assumes that preferences are unconsciously shaped (as in the replicator dynamics) and that human beings are able to perceive altruism and envy from a continuum then one has to admit that using ESS is too fragile. Following these postulates, it has been shown that the strategic environment still determines the evolutionary fitness of altruism and envy. Altruism remains evolutionarily viable if the underlying game exhibits strategic complements and envy remains evolutionarily viable if strategic substitutes are present.

In particular, we have shown that the strict equilbrium preference is CSS and NIS in the original case, and $\mathscr{ER}$ in the Taylor-approximated case with a quadratic fitness function. An interior CSS is an asymptotically stable rest point of the adap-

tive dynamics approach which gives an exact meaning of the evolutionary analysis if restricted to evolution of monomorphic populations. As mentioned before, we are more interested in the consequences for the replicator dynamics where the population composition changes in a distributional sense. In the case that the adopted population space is determined by the strong topology, the dynamic stability of a homogeneous population is given if the corresponding equilibrium strategy is *un-invadable* (which is implied by $\mathscr{ER}$). However, if the weak topology is adopted, the relation to static stability concepts is more restrictive concerning technical and initial state issues. Knowing from Heifetz et al. (2007a) that a somewhat restricted version of the game with full support on the preference space converges under the replicator dynamics, we further can rightly conjecture that the equilibrium preference is dynamically stable.

A few open questions remain. One basic issue concerns the definitions of altruism and envy. Both definitions are fully determined through the subjective well-being function (Eq. (2.11)) (and the sophisticated version of section 2.5) which is somewhat special. One can think of alternative definitions and put these in the discussed evolutionary framework. For example, Heidhues and Riedel (2007) consider *complementary altruism* in the sense of $U_i = \min\{\pi_i, \alpha\pi_j\}$, with some altruism-parameter $\alpha \geq 1$, which follows the ideology of John Rawls in that human beings are rather altruistic if they are better off in relation to others. It would be interesting to see what preferences would emerge in such an alternative game. It would also be of interest to explore whether the evolutionary viability of continuous preferences remains in intermediated cases where the individuals anticipate the opponents' idiosyncratic preferences with a noise term. These issues import chances of future research.

## 2.7 Appendix

*Proof of Proposition 2.4*

To prove $\mathscr{E}\mathscr{R}$ of the equilibrium preference $\alpha^* = \frac{k}{2-k}$ with respect to the relevant fitness function $f(\alpha,\beta) = -\frac{m^2(k(\alpha+1)+2)(k^2\alpha(\beta+1)+k(\alpha-1)-2)}{(4-k^2(\alpha+1)(\beta+1))^2}$, we build on the work of Cressman and Hofbauer (2005) who suggest a classification scheme for testing dynamic stability relating to *quadratic* payoff functions. The basic idea is to consider the payoff function in terms of the mean $E(P)$ and the variance $Var(P)$ of population $P$. In particular, the $k$th moment of $P$ is defined as $P_k = \int x^k P(dx)$ such that $P_1 = E(P)$ and $P_2 = Var(P) + P_1^2$ hold true.

Employing the second-order Taylor expansion of $f(\alpha,\beta)$ around $\alpha = \beta = \alpha^* = \frac{k}{2-k}$ guides to:

$$f(\alpha,\beta) = -\frac{1}{1024(k-1)^3}\left(m^2\left(a\alpha^2 + b\beta^2 + c\alpha\beta + d\alpha + e\beta + f\right)\right) + \text{h.o.t.,} \quad (A.2.1)$$

with coefficients:

$a = -64k^2 - 80k^4 + 128k^3 + 20k^6 - 8k^7 + k^8,$

$b = 48k^4 + k^8 - 12k^7 + 48k^6 - 80k^5,$

$c = 160k^5 + 20k^7 - 160k^4 - 2k^8 - 80k^6 + 64k^3,$

$d = 96k^5 - 128k^4 + 64k^3 - 32k^6 + 4k^7,$

$e = 128k^2 - 256k^3 + 128k^4 - 16k^5 + 8k^6 - 4k^7,$

$f = 256 - 512k + 192k^2 + 64k^3 + 16k^4 - 12k^6.$

For $\mathscr{E}\mathscr{R}$ we have to check whether $f(\delta_{\alpha^*}, Q) = f(\alpha^*, Q) > f(Q,Q)$ holds for all $Q \neq \delta_{\alpha^*}$ sufficiently close to $\delta_{\alpha^*}$ (in the weak topology).

By ignoring the higher order terms and substituting the moments, we reach

$$f(Q,Q) = -\tfrac{m^2}{1024(k-1)^3}\left[aQ_2 + bQ_2 + cQ_1^2 + dQ_1 + eQ_1 + f\right]$$

$$= -\tfrac{m^2}{1024(k-1)^3}\left[(a+b)\,Q_2 + cQ_1^2 + (d+e)\,Q_1 + f\right], \text{ and}$$

$$f(\alpha^*,Q) = -\tfrac{m^2}{1024(k-1)^3}\left[a\,(\alpha^*)^2 + bQ_2 + c\alpha^*Q_1 + d\alpha^* + eQ_1 + f\right],$$

where $-\tfrac{m^2}{1024(k-1)^3} > 0$ by restrictions on $k$.

Thus, the condition of $\mathscr{ER}$ is given by

$$f(\alpha^*,Q) - f(Q,Q) > 0 \Leftrightarrow a\left((\alpha^*)^2 - Q_2\right) + c\,(\alpha^*Q_1 - Q_1^2) + d\,(\alpha^* - Q_1) > 0$$

$$\Leftrightarrow a\left((\alpha^*)^2 - \left(VAR\,(Q) + E\,(Q)^2\right)\right) + c\left(\alpha^*E\,(Q) - E\,(Q)^2\right) + d\,(\alpha^* - E\,(Q)) > 0$$

$$\Leftrightarrow a\left(\left(\tfrac{k}{2-k}\right)^2 - \left(VAR\,(Q) + E\,(Q)^2\right)\right) + c\left(\tfrac{k}{2-k}E\,(Q) - E\,(Q)^2\right) + d\left(\tfrac{k}{2-k} - E\,(Q)\right) >$$

$0$, where we now substitute $\epsilon := VAR\,(Q) > 0$ and $\eta := \tfrac{k}{2-k} - E\,(Q)$, so that

$$a\eta\left(\frac{k}{2-k} + E\,(Q)\right) - a\epsilon + c\eta E\,(Q) + d\eta > 0 \tag{A.2.2}$$

is deciding.

The essence of the weak topology is captured by the fact that populations are close to each other if the respective means are sufficient similar and the variances sufficient close to 0. Hence, for $\mathscr{ER}$, it is left to check whether Ineq. (A.2.2) holds for small $\epsilon > 0$ and small $|\eta|$ where $\eta$ is either positive or negative. Note that we have $a < 0$ with strategical interdependence ($k \neq 0$), such that a sufficient condition becomes

$$\eta\left[a\left(\frac{k}{2-k} + E\,(Q)\right) + cE\,(Q) + d\right] > 0.$$

Accordingly, we have to make a case differentiation regarding $\eta$ and verify that the term $a\left(\tfrac{k}{2-k} + E\,(Q)\right) + cE\,(Q) + d$ is positive for $i)$ $\eta > 0$ and negative for $ii)$ $\eta < 0$.

By inserting the coefficients $a, c, d$ and $E\,(Q) = \tfrac{k}{2-k} - \eta$ into the term, it is straightforward to calculate that

$$a \left( \frac{k}{2-k} + E\left(Q\right) \right) + cE\left(Q\right) + d$$

$$= \underbrace{k^2}_{>0} \left( \underbrace{240\,k^2 - 160\,k^3 + 60\,k^4 - 12\,k^5 - 192\,k + k^6 + 64}_{>0} \right) \underbrace{\eta}_{\substack{i)>0 \\ ii)<0}}.$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{under case } i)\ \ _{>0} \ \ \text{under case } ii)\ \ _{<0}}$$

As a result, Ineq. (A.2.2) is verified, whether the mean of a mutant population $Q$ is slightly below or above the equilibrium. Thus, the Dirac-delta distribution of the homogeneous population $\delta_{\alpha^*}$ with unit mass on the equilibrium preference $\alpha^* = \frac{k}{2-k}$ satisfies the requirements of $\mathscr{ER}$.

$\square$

# Chapter 3

# A Dynamic Model of Reciprocity with Asymmetric Equilibrium Payoffs

*We analyze indirect evolutionary two-player games to identify the dynamic emergence of (strong) reciprocity in a large number of economic settings. The underlying evolutionary environment allows for an arbitrary initial population state provided that every degree of the compact space of reciprocity is adherent to at least one individual of the corresponding continuum population. The basic results, which essentially maintain the evolutionary viability of reciprocity, are, in several directions, context dependent, and minimum valid for the wide class of evolutionary dynamics which hold for regularity and payoff-monotonicity. The evolutionary solution concept which is applied to elevate the explanatory power of emerging Nash equilibria is dominance solvability, in this case, for continuous strategy spaces. An asymmetric aspect comes into play since the actions of the evolutionary players are not only determined by the current state of reciprocity but also by their inherent, context-free preferences towards others which differ among one another devoid of being endogenized in the time span of the dynamic process at hand.*

## 3.1 Introduction

In recent years, it has been established that the orthodox assumption on material or monetary[1] selfishness of players is not necessarily sustainable in economic modeling. While the assumption of exogeneous given selfishness fits fairly well in some economic contexts,[2] real life evidence and experimental data suggest that people do not behave consistent with this postulate in general. Most convincing studies include ultimatum/dictator games (e.g. Güth et al., 1982; Andreoni and Miller, 2002; Camerer, 2003), or public goods contribution games (e.g. Andreoni et al., 2002).

The present research contributes to the literature by analyzing a dynamic model of reciprocity in an evolutionary framework.[3] In essence, reciprocity refers to peoples' desire to reward perceived kindness and to punish perceived unkindness. The present form of reciprocity invokes the idea that subjective values include to be sensitive to opponents' intrinsic preferences. More precisely, our suggestion of players' *psychological payoff* including other-regarding motivations follows the model of Levine (1998) who demonstrates a striking consistency with results from experimental lab studies. Levine shows that his model is useful to understand several results from ultimatum games and market experiments. Formally, player $i$ seeks to

---

[1]Putting the meaning of selfishness to an economic environment is essential for the present study and related models with so-called *other-regarding* preferences. To make this point, consider Joel Sobel's comment on the hypothesis *only the selfish survive*: "with sufficient freedom to define "selfish" this statement is a tautology" (Sobel (2005, p. 430)). In other words, motivations may depend on others' motivations but their exclusive personality is a banality in the final analysis. This thinking appears trivial yet corresponds to a famous theory called "psychological egoism" which claims that anything we do for others is just because of increasing our own welfare.

[2]Generally, in highly competitive settings with many players and with one-shot and/or anonymous play (e.g. financial markets), one can assume that the average behavior reflects a high degree of material greed.

[3]As described in the next section, one should be aware of the fact that reciprocity has different definitions in the economic literature. See also the survey paper by Sobel (2005) for this issue.

maximize her subjective well-being, given by[4]

$$v_i = u_i + \sum_{j \neq i} \frac{\alpha_i + \lambda \cdot \alpha_j}{1 + \lambda} \cdot u_j,$$

where $u_i$ and $u_j$ are material payoffs, $\alpha_i, \alpha_j \in (-1, 1)$ are social preference parameters, and $\lambda \in [0, 1]$ symbolizes a weight which player $i$ puts on player $j$'s preference. Accordingly, individuals are not only concerned with their own material payoff but also with that of their opponents. This fact is in the first place due to intrinsic preferences like altruism. However, by considering an extra dimension of reciprocity[5], the individual weight which is placed on the opponents' profit varies additionally with respect to the opponents' intrinsic preference. The advantage of the well-being functional with *two* preference dimensions is that it allows us to explore why people sometimes behave contrary to their true attitudes or *ethos*. In particular, the daily observance, whether in economic or other social life settings, of intrinsic good people behaving badly (or selfishly) sometimes or intrinsic selfish (or bad) people behaving well sometimes centers the motivation of the current study.

Technically, we use an indirect evolutionary environment to identify the cultural viability of reciprocity in a broad class of pre-programmed populations where the players engage in many different types of strategic two-player games. Indirect evolution allows the players to choose a behavior which follows their perceptive payoffs but the players receive evolutionary fitness (or reproductive success) according to their "true", objective payoffs. The presupposition that players aim to maximize their own idiosyncratic preferences but that economic success "regulates the market" is by now a central tenet in economic modeling and successfully opposes the

---

[4]This functional is the exact writing of Levine (1998, p. 597). Sethi and Somanthan (2001) introduce a similar model by replacing player $i$'s weight on player $j$'s material profit $[\alpha_i + \lambda \cdot \alpha_j] / [1 + \lambda]$ by $[\alpha_i + \lambda \cdot (\alpha_j - \alpha_i)] / [1 + \lambda]$, where $0 \leq \alpha_i < 1$ and $\lambda \geq 0$. Their specification allows an altruist to place a negative weight on a selfish individual which is not possible under Levine's definition. Apart from the denotative conception of reciprocity and for the sake of technical simplicity, the precise definition of the players' subjective payoffs are once more slightly modified in the present study.

[5]Although Levine does not explicitly connect with the term "reciprocity" in his study (however, Sethi and Somanthan do), the parameter $\lambda$ clearly represents much of the features of reciprocity in any environment of non-anonymous interaction.

neo-classical feature of pure economic selfishness. Consequently, a reciprocity disposition which yields a larger objective payoff tends to become more prevalent in a certain population while dispositions with relative low objective payoffs tend to decease. The idea that social preferences beside selfishness are profitable in some strategic environments is well documented in the relevant literature and dates back at least as far as Schelling (1960). The basic reason for this is that (at least somewhat) recognized preferences provide a commitment device which makes unselfish players potentially the better performers in interdependent settings. Another precursor of this theme is Frank (1987, 1988) who finds that recognizable emotions have the strategic power to change the actions of others, and therefore features the ability to increase the profit of the possessors. Güth and Yaari (1992) then formally introduce the approach of indirect evolution. In fact, they use this approach to challenge a form of reciprocity by having regard for individuals who have tendencies for rejecting unfair offers in ultimatum bargaining. However, the reciprocity motive is somewhat restricted there because the agents are not able to feel subjective benefits from proposing fair distributions. In line with related literature, the finding of Güth and Yaari is that the observability of types guarantees that reciprocators gain from their attitudes in material terms while opportunists relatively loose.[6]

In the present study, we analyze the evolving of reciprocity in the wide frame of regular and payoff-monotonic selection processes. For that purpose, we use a result which is shown by Heifetz et al. (2007a): if the evolutionary game on the level of biases (which is located at equilibrium behavior which results from the players' idiosyncratic well-being functionals) is dominance solvable, in the sense of Moulin (1984), for continuous preference spaces, then, the limiting population can be characterized under any payoff-monotonic selection dynamics. The replicator dynamics for general distributions (cf. Oechssler and Riedel (2001) for a detailed technical

---

[6]Huck and Oechssler (1999) illustrate viability of preferences for rejecting unfair divisions even if preferences are unobservable provided that the population is small and that the distribution of preferences is known in the population.

survey) is subsumed by the class of payoff-monotonic growth-rate functions that is applied to the present study. The more general class considered here explicitly allows to interpret the evolving of reciprocity not just on a biological level but rather on a cultural level of learning (including imitation) or education where a reciprocity disposition that is adequate to achieve higher material returns is replicated faster. Note that our approaching is somewhat enhanced in comparison with the practice of a group of research articles that assumes the evolving of other-regarding issues (like reputation, social preferences, positional goods, ideologies, and so on) on a biological level and treats the results simply as a metaphor to interpret the dynamics in terms of cultural spreading.[7] Yet, this is not the only reason why the present methodological approach is prepared for exploring reciprocity. In alignment with the dynamical system, we are allowed to assume *any* arbitrary initial population distribution provided that it is described by a compact interval of reciprocity. This technical feature equips the results with a striking general character.

A first observation of our work shows that a sufficient reciprocal player who is intrinsically spiteful (altruistic) may behave benevolent (spiteful) if the opponent player is sufficient contrary in the intrinsic attitude and the strategic situation requires that kind of behavior. In fact, the sign of the players' overall concern for the other players' payoff is always determined by the strategic situation, i.e. the players show a positive concern for the other players if strategic complements are present and a negative concern if strategies are substitutes—a result which is reminiscent to related work of Possajennikov (2000) and Bester and Güth (1998). A further observation says that if player A's intrinsic preference lies below a certain threshold

---

[7]Answering the question whether preference evolution relies on a biological or cultural selection process is a subtle task and rarely elucidated in much detail in related work. However, in our case, it is indeed useful to examine reciprocity on a rather short-termed cultural level since we assume two differently treated preference dimensions which initiate the equilibrium actions: one (altruism/spite) is exogenous given and not evolving while the other (reciprocity) is endogenized and evolving. With these specifications, it is appropriate to think of reciprocity as a cultural norm which changes within the time span of cultural spreading while altruism/spite changes more seldom by gene transmission. For a deeper understanding of the different selection processes (and the associated speed differences) in evolutionary game theory, see Selten (1991).

(if player A is relatively spiteful), then player A's tendency to reciprocity is higher the lower the intrinsic preference of player B (the nastier player B)—if player A's intrinsic preference lies above the threshold, then player A's tendency to reciprocity is higher the higher player B's intrinsic preference is. The threshold depends on the strategic type of the underlying game.

The remainder of the paper works as follows. The first part of section 3.2 surveys some recent approaches of reciprocity or fairness in economics which are, in some obvious sense, connected with the present study. The second part illustrates why our treatment of subjective payoff perception, involving reciprocity, has the strategic power to resolve social dilemma conflicts. Section 3.3 then introduces the basic model under which the viability of reciprocity is analyzed. Section 3.4 concludes. Technical details about the applied dynamics, and a figure and a table, that specify some initial conditions, appear in the appendices.

## 3.2 On Reciprocity

### 3.2.1 Related Models and Current Treatment

In order to explain the emergence of other-regarding preferences, there are by now several studies that try to identify more or less complicated models which give insight into the economic psychology of people in strategic situations. The common theme in these models is the antithesis to the neo-classical assumption that people's behavior is thoroughly driven by material selfishness. In order to narrow the wide spectrum of recent approaches and to connect with the present work, it is useful to concentrate on prominent models which explicitly incorporate a certain motive of reciprocity or *fairness*.[8] One such class assumes that subjective benefits

---

[8]Usually, the will to reciprocate springs from the will to being fair. However, the concept of fairness is likewise of somewhat ambiguous use in economics. At this point, we refrain from a broader discussion on fairness and point to the well-being functionals of this section for examining the specific conception.

are motivated from inequity aversions of own and other's economic gains. Put differently, the players' actions are initiated by distribution considerations. Fehr and Schmidt (1999) propose for this approach. In their regard, the subjective well-being functional of player $i$ in the standard two-player setting is formalized as

$$U_i = \pi_i - \alpha_i \max[\pi_j - \pi_i, 0] - \beta_i \max[\pi_i - \pi_j, 0],$$

where $\pi_i, \pi_j$ are material payoffs and $\alpha_i \geq \beta_i \geq 0$, $\beta_i < 1$ are weight parameters. Hence, the perceptive payoff of player $i$ differs significantly in the issue whether she and her opponent $j$ are approximately equal rich in economic values; viz., the players are pre-programmed to feel satisfaction from being about as rich as the opponent players. To further summarize, the players are inequity-averse ($\alpha_i, \beta_i \geq 0$), dislike inequity more if it springs from own relative loss ($\alpha_i \geq \beta_i$), but like gaining profit more than reducing inequity ($\beta_i < 1$). A similar model, also motivated by the idea that the players aim to reduce inequity in their material payoffs, is developed by Bolton and Ockenfels (2000). In the two-player setting, Bolton and Ockenfels assume that the personal well-being of player $i$ is determined via the (possibly non-linear) term

$$U_i = v_i \left( \pi_i, \frac{\pi_i}{\pi_i + \pi_j} \right),$$

where $v_i(\cdot, \cdot)$ is globally non-decreasing, concave in the first argument (the material payoff of player $i$), and strictly concave in the second argument (the relative material payoff of player $i$). The models of Fehr/Schmidt and Bolton/Ockenfels are both motivated by the idea that players act according to satisfy their fairness emotions by reducing economic inequity. However, the players in these models are distributional motivated and do not explicitly estimate the individual types of the opponents, i.e. the players do not differentiate in the other players' intentions or preferences. Inspired by the psychological game-theoretical approach of Geanakoplos et al. (1989), there is a somewhat more complex class of fairness models which accounts for these elements and seems to reflect reality more detailed. In psychological games, the players' preferences depend on their beliefs about the other players'

52

intentions. By using normal form games, Rabin (1993) proposes for this technique. With his notation, he assumes that individual $i$ plays according to her expected utility

$$U_i\left(a_i, b_j, c_i\right) = \pi_i\left(a_i, b_j\right) + \widetilde{f}_j\left(b_j, c_i\right)\left[1 + f_i\left(a_i, b_j\right)\right],$$

where $a_i, b_j,$ and $c_i$ are, in this order, the strategy chosen by player $i$, the belief of player $i$ about the strategy of player $j$, and the belief of player $i$ about the belief of player $j$ about the strategy of player $i$. $\pi_i\left(a_i, b_j\right)$ is player $i$'s material payoff, $\widetilde{f}_j\left(b_j, c_i\right)$ is player $i$'s belief about the kindness of the opponent $j$ towards player $i$, and $f_i\left(a_i, b_j\right)$ symbolizes the kindness of player $i$ towards player $j$. If equilibrium play is reached, the players' beliefs about the other players' intentions are true and the players base their actions on these beliefs and the subsequent actions of the other players. In line with Rabin's approach but with the purpose to expand to extensive form games, Dufwenberg and Kirchsteiger (2004) assumes a similar model. The basic difficulty of the extensive form is given by the fact that the players have to adjust their perceived utility by updating their beliefs about the others' intentions at each node of the sequential game tree to ensure useful results. Contrary, in Geanakoplos et al. (1989) and Rabin (1993), the players have only initial beliefs about the others' intentions which make the equilibrium analysis easier. Falk and Fischbacher (2006) and Charness and Rabin (2002) propose equilibrium models which basically incorporate both aspects the distributional one and the intention-based. However, opposing to the approach used in the current paper, the applicability of these models is somewhat limited which is basically due to the assumption of higher order beliefs about the others' intentions and the appearance of many equilibria. Levine's model and the version used here depict a third way of modeling reciprocity. In particular, the subjective utility functions of the players depend on the beliefs about the intrinsic preferences of others, i.e. the players reciprocate to the perceived preference of the respective opponent player.

From the previous sentences, it becomes clear that the meaning of reciprocity

is ambiguous in terms of functional forms subsuming a motive of fairness. However, discriminations of reciprocity are multi-dimensional existent. In the following, we will mention some further basic facets which appear repeatedly throughout the literature. One prominent aspect is to distinguish between weak and strong reciprocity. Weak reciprocity stands for the conception that people reciprocate in order to gain higher material returns in the future by sustaining collaboration. Typically, weak reciprocity relies on reputation and repeated interaction in orthodox economic modeling and is basically not different from pure selfishness in social preference terminology. In a key paper, Trivers (1971) uses the term *reciprocal altruism* which is identical to weak reciprocity for which he shows sustainability under infinite repeated interactions.[9] In contrast, strong reciprocity refers to the conception that people show cooperative or retaliatory behavior, even if there is no reason to expect higher material returns in the future (e.g. Gintis, 2000). Moreover, people are willing to sacrifice own profit in order to either help friends or harm enemies. Under this aspect, strong reciprocity is *really* other-regardingly intended.[10]

The differentiation of positive and negative reciprocity among the literature is self-explanatory (e.g. Hoffmann et al., 1998). Loosely speaking, positive reciprocity describes the tendency to reward kind people while negative reciprocity describes the tendency to harm cruel people.

Another common aspect is to differentiate between direct and indirect reciprocity (e.g. Nowak and Sigmund, 2005). Direct reciprocity describes the routine that if 'person A helps (harms) person B, then person B helps (harms) person A' while indirect reciprocity states that if 'person A helps (harms) person B, then person C

---

[9]As already discussed in the literature, the denotation "reciprocal altruism" appears somewhat inadequate in this respect, cf. Hoffmann et al. (1998, p. 338). They argue convincingly that "I am not altruistic if my action is based on my expectation of your reciprocation". Another thread of research comments that weak reciprocity is not reciprocity and would therefore probably be unhappy with Trivers' denotation even in this aspect. For example, Fehr and Fischbacher (2002, C3) write: "It is important to emphasize that reciprocity is not driven by the expectation of future material benefit. It is, therefore, fundamentally different from "cooperative" or "retaliatory" behavior in repeated interactions."

[10]For obvious reasons, Sobel (2005) substitutes "strong" with "intrinsic", and "weak" with "instrumental".

helps (harms) person A'.

The specific notion of reciprocity used in the current paper is determined by the subjective well-being functionals of section 3.3 and the underlying methodological approach, and, in this regard, described as follows.

In comparison with the significant recurrent features of economic reciprocity, the present shape exhibits the following attributes:

- *preference-based*

- *strong* (*intrinsic*)

- both, *positive or negative*

- *indirect.*

Though, the stated attributes are not intended to identify an objective, "true" definition of reciprocity in economics. More precisely: the aim of this paper is not to elucidate the meaning of reciprocity how it should be used in economics but rather to assume a form of reciprocity that exhibits some predominant features which appear frequently in the literature, and to explore under what circumstances this form of reciprocity can survive in an evolutionary process.

### 3.2.2   A Simple Illustration

As we will see, the implications of strong reciprocity which we obtain in our basic model are somewhat intricate to retrace (however, the trend and interpretation of these results remain on a plain level). So, the following formulation of the symmetric two-player prisoners' dilemma is intended to give a simple illustration why the current treatment of subjective payoff perception, involving reciprocity, has the strategic power to resolve social dilemma conflicts and overcome spite.[11] Consider

---

[11]The prisoners' dilemma is the leitmotif in Sethi and Somanthan (2003) for surveying economic reciprocity in the evolutionary game theoretic literature. The current version uses a different notion of reciprocity.

the following $2 \times 2$ matrix.

|            | cooperate              | defect       |
|------------|------------------------|--------------|
| cooperate  | $\xi - \upsilon, \xi - \upsilon$ | $-\upsilon, \xi$ |
| defect     | $\xi, -\upsilon$       | $0, 0$       |

Figure 3.1: A Prisoners' Dilemma

As is the rule in the matrix design, one player is the row player and the other plays the column; the first number in each matrix entry is the payoff received by the row player and the second one belongs to the column player. The strategy *cooperate* is connected with a private loss of $\upsilon > 0$ and a benefit to the other player of $\xi > \upsilon$. The strategy *defect* yields neither a loss nor a benefit. If two intelligent and self-interested individuals play this game exactly once, we face the well-known dilemma where both players *defect* in order to reach a higher payoff regardless of the strategical choice of the other player. However, *cooperate* would be mutually better since both players' outcome is higher under this strategy profile: $\xi - \upsilon > 0$ with $\xi > \upsilon$.

Under a simple model of natural selection the same dilemma defines the usual situation. If we think of a population with size $N$ consisting of $k$ cooperators, and hence $N - k$ defectors, the reproductive matrix payoffs (or *fitnesses*) are given by

$$\mathfrak{f}_C = \frac{k-1}{N-1}\xi - \upsilon \quad \text{(cooperators)}$$

and

$$\mathfrak{f}_D = \frac{k}{N-1}\xi \quad \text{(defectors)},$$

and the average fitness is determined by $\bar{\mathfrak{f}} = \frac{k}{N}(\xi - \upsilon)$. In any mixed population the defectors reach a higher fitness than the cooperators so that natural selection tends to decline the fraction of cooperators while the fraction of defectors eventually take over the whole population. Hence, without any model arrangement which favors the outcome of cooperation, the dilemma of the one-shot game trivially persist under natural selection where matrix payoffs correspond to fitnesses.

Consider now the notion of player $i$'s perceived payoff which we use in the following chapters, i.e.[12]

$$U_i = \pi_i + \theta_i \cdot \pi_j \text{ with } \theta_i = \alpha_i \cdot \gamma_i + (1 - \alpha_i) \cdot \gamma_j,$$

where $\alpha_i \in [0, 1]$ defines the individual "norm of reciprocity", and $\gamma_i, \gamma_j \in [-1, 1]$ defines the individual "intrinsic preference".[13] Note that if $\gamma_i \neq 0$, then player $i$ puts a non-zero weight on the material payoff of player $j$ (unless the extremly rare case where $\alpha_i \cdot \gamma_i + (1 - \alpha_i) \cdot \gamma_j = 0$ with $\gamma_i \neq 0$; however, even in this case individual $i$ is intrinsically biased but her disposition does not come into effect only because her reciprocity norm and the intrinsic preferences of both players compensate to zero). Hence, we say that player $i$ is biased if $\gamma_i \neq 0$. Further, we say that player $i$ is materialistic if both hold true $\gamma_i = 0$ and $\alpha_i = 1$, since then $\theta_i = 0$, i.e. the material outcome of player $i$ coincides with her perceived payoff. Accordingly, a materialist places no weight on the other players' payoff while a biased player $i$ places a weight of $\rho_i^b$ on the payoff of a biased player and a payoff of $\rho_i^m$ on the payoff of a materialist, where

$$\rho_i^b = \alpha_i \cdot \gamma_i + (1 - \alpha_i) \cdot \gamma_j; \quad \rho_i^m = \alpha_i \cdot \gamma_i.$$

If two biased player interact, we reach the following payoff matrix.

|  | cooperate | defect |
|---|---|---|
| cooperate | $\xi - \upsilon + \rho_i^b (\xi - \upsilon), \xi - \upsilon + \rho_j^b (\xi - \upsilon)$ | $-\upsilon + \rho_i^b \xi, \xi - \rho_j^b \upsilon$ |
| defect | $\xi - \rho_i^b \upsilon, -\upsilon + \rho_j^b \xi$ | $0, 0$ |

Figure 3.2: A Prisoners' Dilemma with Biased Players

Provided that $\xi - \upsilon + \rho_\bullet^b (\xi - \upsilon) > \xi - \rho_\bullet^b \upsilon$ and $-\upsilon + \rho_\bullet^b \xi > 0$, and thus $\rho_\bullet^b > \frac{\upsilon}{\xi}$, *cooperate* is a dominant strategy for both players (the $\bullet$ stands for either $i$ or $j$). Note that if ① $\alpha_i \cdot \gamma_i + \gamma_j > \alpha_i \cdot \gamma_j$ and ② $\alpha_j \cdot \gamma_j + \gamma_i > \alpha_j \cdot \gamma_i$ then $\rho_\bullet^b > 0$ and

---

[12]Cf. Eqs. (3.1) and Eqs. (3.3) in section 3.3. For the sake of simplicity, we assume that all dispositions are perfectly observable.

[13]In our main model in section 3.3 we exclude *perfect* intrinsic preferences for technical reasons. This is not necessary for the current illustrative purpose.

($cooperate, cooperate$) is potentially a Nash equilbrium, depending on the ratio of benefit and loss. It is easily comprehended that both inequalities ① and ② hold if both players have altruistic feelings towards others, i.e. $\gamma_i, \gamma_j > 0$. But even in the case of different intrinsic preferences, i.e. $\text{sign}(\gamma_i) \neq \text{sign}(\gamma_j)$, the reciprocity motive is apt to keep *cooperate* as the agreed strategy. For example, if player $j$ is moderate spiteful, say $\gamma_j = -0, 6$, and player $i$ is perfectly altruistic, $\gamma_i = 1$, a high reciprocity norm of player $j$, say $\alpha_j = 0, 3$, can reverse player $j$'s natural will to defect. Note that the ability to reciprocity gives a strong impetus to the game. Even if a player has a strong intrinisic attitude ($|\gamma_i|$ is close to 1) a perfect tendency to reciprocity ($\alpha_i = 0$) will always overcome the origin will to either cooperate or defect if the other player has a contrary intrinsic attitude ($\text{sign}(\gamma_i) \neq \text{sign}(\gamma_j)$) since then $\text{sign}(\theta_i) = \text{sign}(\gamma_j) \neq \text{sign}(\gamma_i)$. Of course, if both players are spiteful, i.e. $\gamma_i, \gamma_j < 0$, then ① and ② never hold and the only rational strategy is always defect.

In the case that a biased player $i$ meets a materialistic player $j$, the action of the biased player is determined by her intrinsic preference and independent of her reciprocity motive, since the sign of $\rho_i^m$ (and thus the sign of $\gamma_i$) induce whether to cooperate or defect. Naturally, a materialist always defects and is in the advantageous free-reding position if the opponent is a cooperating benevolent player. Clearly, the pros and cons of being biased in the matrix PD game continue in a standard population model where the matrix payoffs correspond to fitnesses.

The prisoners' dilemma essentially illustrates the strategic advantages which can result from the reciprocity motive when two biased players interact where one player is sufficient altruistic and the other is intrinsically malevolent but sufficient reciprocal to overcome this attitude. It also demonstrates that two altruists can always resolve the social dilemma but the reciprocity motive is not apt to overcome the will of an altruist to cooperate if the other player is materialistic in the above sense. Hence, the expected evolutionary advantages which results from the reciprocity motive seem to depend heavily on the initial population distribution regarding the intrinsic preferences of the players.

## 3.3  Model

In this section we will introduce our model of strong reciprocity, state our main result under the assumption that the players perfectly recognize each others types, and interpret the results.

Let there be a large population of evolutionary agents. At each instant in time a pair of agents is matched at random to play the game $\Gamma_U = (\{1,2\}, \{x,y\}, \{U_1, U_2\})$ with the aim to maximize their subjective well-being, determined by

$$U_1 = \pi_1 + \theta_1 \cdot \pi_2 \tag{3.1a}$$

$$U_2 = \pi_2 + \theta_2 \cdot \pi_1, \tag{3.1b}$$

where $\pi_1$ and $\pi_2$ are material or "economic" (and therefore interpersonal comparable) payoffs (e.g. money). In order to incorporate a broad variety of strategic situations in this study, we assume that the material payoffs are defined by

$$\pi_1 = x \cdot (l \cdot y - x) + x \tag{3.2a}$$

$$\pi_2 = y \cdot (l \cdot x - y) + y, \tag{3.2b}$$

where $x, y \in [0, \infty)$ describe the actions or efforts of player 1 and player 2, respectively. The parameter $l \in \left(\underline{l} < 0, \bar{l} > 0\right)$ determines the characteristic nature of the game by measuring the kind and extent of strategical interdependence; the $l$ is further specified as soon as required. The specification of the economic payoffs is sufficiently general to illustrate success since the economic interpretations are extensive. The simplest example is a production game with either negative or positive externalities which is determined by the sign of $l$. The externality $l < 0$ represents, for example, a common pool resource game where the players exploit a resource with efforts $x$ and $y$, respectively. Accordingly, the higher player A's input the lower player B's payoff. Oppositely, $l > 0$ determines a game where a more aggressive behavior of one player increases the payoff of the other player like in public good

contribution settings. Alternatively, one can assume oligopolistic competition where the efforts are either firms' quantity choices (in a Cournot market) or price choices (in a Bertrand market).

The variables $\theta_1, \theta_2$ symbolize the subjective overall concern for the respective opponents' profit, and have deeper meanings, as specified by

$$\theta_1 = \alpha \cdot \gamma_1 + (1 - \alpha) \cdot \gamma_2 \tag{3.3a}$$

$$\theta_2 = \beta \cdot \gamma_2 + (1 - \beta) \cdot \gamma_1, \tag{3.3b}$$

where the parameters $\gamma_1, \gamma_2$ are intrinsic preferences or attitudes, either altruism or spite[14] ($\gamma_1, \gamma_2$ include also material selfishness at the peak of neutrality). The variables $\alpha, \beta$ identify the dispositions to reciprocity which belong to player 1 and player 2, respectively. While the dimension of altruism and spite is an exogenous trait, the dimension of reciprocity is endogenized in the model. This distinction allows, for example, a player who is rather altruistic inclined to behave spiteful in a reciprocal manner, however, by keeping the true character. Or, a player who is rather spiteful by nature is able to behave benevolent without changing the true character during the course of selection. From these specifications, one should think about reciprocity as a cultural norm or convention which is changeable by the dynamical pressures, and in this line, provides the distinct flexibility in the players' behavior. It is more for the sake of distinctiveness that we will sometimes refer to the parameters $\gamma_1, \gamma_2$ as intrinsic preferences and to the reciprocity-variables $\alpha, \beta$ as cultural norms or conventions; because, in the sense of subjective motivations which distort the economic greed of the players and initiate their actions, $\alpha$ and $\beta$ belong to the class of intrinsic preferences, too. This is rather a question of definition.

In accordance with Levine (1998) and Sethi and Somanthan (2001), we avoid initially the somewhat unnatural economic situations in which a player is less (or equally) concerned about herself than about the opponent. Formally, the subjective

---

[14]Alternatively, one can assume that envy or malevolence is the opposite preference to altruism.

overall concern for the other player satisfies $|\theta_1|, |\theta_2| < 1$. To ensure this, we impose the following restrictions on the preference components:

$$\gamma_1, \gamma_2 \in A = [-1 + \epsilon, 1 - \epsilon]$$
$$\alpha, \beta \in B = [0, 1], \tag{3.4}$$

where $\epsilon$ is positive and small. Both assumptions are intuitively plausible. The first one allows the players to exhibit a negative ("spite": $\gamma_1, \gamma_2 < 0$), a neutral ("egoism": $\gamma_1, \gamma_2 = 0$), or a positive ("altruism": $\gamma_1, \gamma_2 > 0$) intrinsic preference towards others. The second assumption is even more intuitive and shows that the players evaluate their overall concern by a convex combination of the own and the other players' intrinsic preference. Note that $\alpha$, $(\beta)$ induce reciprocal actions only if $\alpha$, $(\beta) \neq 1$ since in the case of $\alpha$, $(\beta) = 1$ the agents' subjective overall concern is independent of the opponents' intrinsic preference. Accordingly, we have the following notion.

**Definition 3.1.** *The agents possess a tendency to reciprocity whenever* $\alpha$, $(\beta) \in B \setminus \{1\}$.

Evidently, the intuition of these assumptions is in line with the restriction of the subjective overall concern towards others, which is finally fixed by $|\alpha \cdot \gamma_1 + (1 - \alpha) \cdot \gamma_2| = |\theta_1|$, $(|\beta \cdot \gamma_2 + (1 - \beta) \cdot \gamma_1| = |\theta_2|) < 1$. More precisely, with intrinsic traits of altruism and spite and the present distribution of reciprocity, the players are completely identified over the compact space

$$\theta_1, \theta_2 \in \Theta = [-1 + \epsilon, 1 - \epsilon]. \tag{3.5}$$

The game setup is close to the one of Harrison and Villena (2008) but differs significantly in several aspects. First, Harrison and Villena concentrate on game settings which exhibit negative externalities ($\frac{\partial \pi_1(x,y)}{\partial y} < 0, \frac{\partial \pi_2(x,y)}{\partial x} < 0$) and strategic substitutes ($\frac{\partial \pi_1(x,y)}{\partial x \partial y} < 0, \frac{\partial \pi_2(x,y)}{\partial x \partial y} < 0$). This means that a higher input of player A lowers both the actual payoff and the marginal payoff of player B. From Eqs. (3.2) and their first and second order derivatives it is easy to see that the present setting

61

represents negative externalities ($\frac{\partial \pi_1(x,y)}{\partial y} < 0$, $\frac{\partial \pi_2(x,y)}{\partial x} < 0$) and strategic substitutes ($\frac{\partial \pi_1(x,y)}{\partial x \partial y} < 0$, $\frac{\partial \pi_2(x,y)}{\partial x \partial y} < 0$) if $l < 0$, and positive externalities ($\frac{\partial \pi_1(x,y)}{\partial y} > 0$, $\frac{\partial \pi_2(x,y)}{\partial x} > 0$) and strategic complements ($\frac{\partial \pi_1(x,y)}{\partial x \partial y} > 0$, $\frac{\partial \pi_2(x,y)}{\partial x \partial y} > 0$) if $l > 0$.[15] Note that with $l = 0$ there is no strategic interdependence so that economic competition becomes "monopolistic". Consequently, the present model incorporates a much broader class of strategic games.

A second difference regards the evolutionary analysis. While Harrison and Villena use the ESS concept to illustrate the evolutionary viability of reciprocity, we use dominance solvability as proposed by Heifetz et al. (2007a). The lack of ESS is that its predictions are only static. ESS does not explore to what level evolution will lead the evolving trait of a certain population but can only tell whether a somehow reached population state is immune to rare "mutations".[16,17] Contrary, dominance solvability is useful to establish dynamic results with respect to many initial population states that evolve according to the broad class of regular, payoff-monotonic dynamics.

According to the indirect evolutionary approach, the players maximize their subjective well-being which leads to a second stage game located at equilibrium behavior (or *on the level of biases*). Let this game be symbolized by $\Gamma_f = (\{1,2\}, \{\alpha, \beta\}, \{f_1, f_2\})$, where $f_1, f_2$ are the reproductive success defining fitness functions that can be identified by substituting equilibrium behavior in the economic payoffs (Eqs. (3.2)), which formally corresponds to

$$f_1, f_2 = \pi_1(x^*, y^*), \pi_2(x^*, y^*),$$

where $x^*, y^*$ are equilibrium strategies of $\Gamma_U = (\{1,2\}, \{x,y\}, \{U_1, U_2\})$. Thus, let

---

[15]The terminology to characterize the strategic environment was introduced by Bulow et al. (1985) in order to distinguish games with upward sloping best-response functions from those with downward sloping best-response functions.

[16]Formally, a strategy $x^*$ is ESS, if either i) $\pi(x^*, x^*) > \pi(x, x^*)$ or ii) $\pi(x^*, x^*) = \pi(x, x^*)$ and $\pi(x, x) < \pi(x^*, x)$ for all mutations $x \neq x^*$, see Maynard-Smith and Price (1973).

[17]Another lack of ESS is the insufficiency for characterizing dynamic stability of certain evolutionary dynamics like replicator or BNN with continuous strategy sets (cf. Hofbauer et al., 2009, and some references therein).

both players maximize their perceived payoffs, i.e. $x^* \in \text{argmax}_x \, U_1(x, y^*)$ and $y^* \in \text{argmax}_y \, U_2(x^*, y)$, which defines their reaction functions: $x = \frac{1}{2}(1 + ly(1 + \theta_1))$ and $y = \frac{1}{2}(1 + lx(1 + \theta_2))$. Equalizing the reaction functions identifies the unique equilibrium profile of the game $(x^*, y^*)$, where

$$x^* = -\frac{\theta_1 l + l + 2}{l^2 + \theta_2 l^2 - 4 + \theta_1 l^2 + \theta_1 l^2 \theta_2} \tag{3.6a}$$

$$y^* = -\frac{\theta_2 l + l + 2}{l^2 + \theta_2 l^2 - 4 + \theta_1 l^2 + \theta_1 l^2 \theta_2}. \tag{3.6b}$$

From the equilibrium profile, we can comprehend the strategic influence of player A's regard for player B's payoff on player B's strategy. The strategic influence is consistent with the psychological idea that individuals condition their actions on the perceived types of others and do not act uniformly with each other. At this point, it becomes clear that the relatedness of the other players' type and the own equilibrium action requires a positive degree of recognition. Note again that we have assumed this ability of the players in the perfect sense.

Plugging the equilibrium actions in the material payoff functions leads to the individual fitnesses which are functions of the biases,

$$f_1\left(\theta_1\left(\gamma_1, \gamma_2, \alpha\right), \theta_2\left(\gamma_1, \gamma_2, \beta\right)\right) = \pi_1(x^*, y^*)$$
$$= -\frac{(\theta_1 l + l + 2)(-l + \theta_1 l - 2 + \theta_1 l^2 + \theta_1 l^2 \theta_2)}{(l^2 + \theta_2 l^2 - 4 + \theta_1 l^2 + \theta_1 l^2 \theta_2)^2} \tag{3.7a}$$

$$f_2\left(\theta_1\left(\gamma_1, \gamma_2, \alpha\right), \theta_2\left(\gamma_1, \gamma_2, \beta\right)\right) = \pi_2(x^*, y^*)$$
$$= -\frac{(\theta_2 l + l + 2)(-l + \theta_2 l - 2 + \theta_2 l^2 + \theta_1 l^2 \theta_2)}{(l^2 + \theta_2 l^2 - 4 + \theta_1 l^2 + \theta_1 l^2 \theta_2)^2}. \tag{3.7b}$$

Eqs. (3.7), the *equilibrium payoffs*, are the central functionals which measure the prevalence of the different types in the game (the specifications of the dynamic process—where successful types proliferate at the expense of abortive types—are given in the Appendix A).

In the following, we assume that the intrinsic preference of player 1 is not exactly the same as the intrinsic one of player 2, i.e. $\gamma_1 \neq \gamma_2$. This assumption gives the

game $\Gamma_f$ an asymmetric character and is reasonable in order to adopt reciprocity in the model. In the case of $\gamma_1 = \gamma_2$, it would be sufficient to behave according to the intrinsic preference altruism or spite to fulfill the characteristic of reciprocity. Formally, there are now two different populations but with intrinsic traits of altruism and spite selected from the same pool. Somewhat informal, one can imagine that nature picks $\gamma_1, \gamma_2$ from the equal distributed set $A$ "with two hands at once". Based on the usual asymmetric setting in evolutionary games (cf. Selten, 1980; Weibull, 1995, pp. 64), let us imagine an ex ante symmetric game, denote $\Gamma_{\Gamma_f}^\gamma$. In this game any intrinsic preference parameter is assigned to each of the players with the same probability. This assumption corresponds to "nature plays first" by allocating $\gamma_1, \gamma_2$ to player 1 and player 2. Relying on Selten's work, the pair of reciprocity biases $(\alpha, \beta)$, where $\alpha$ is associated with $\gamma_1$ (i.e. the biases of player 1 are given with $\gamma_1$ and $\alpha$) and $\beta$ is associated with $\gamma_2$, would be evolutionarily stable in the sense of ESS in the ex ante symmetric game $\Gamma_{\Gamma_f}^\gamma$ if and only if the vector $(\alpha, \beta)$ describes a strict Nash equilibrium of the asymmetric game $\Gamma_f$. However, as mentioned before, the question of interest regards the conception of dominance solvability, and hence, the question of which type pass the dynamic evolutionary pressures under many starting conditions.

The following lemma is useful to identify a dominance solvable trait (cf. Heifetz et al., 2007a; Moulin, 1984, Theorem 4).

**Lemma 3.1.** In order to check for dominance solvability of a particular trait it is sufficient to compute that

  (i) the fitness function is continuous, twice differentiable and strictly concave in the particular trait of each player;

 (ii) the slope of each player's best-reply function is less than 1 in absolute value;

and to argue that

(iii) the particular trait is selected from a compact interval.

To start with, a substantial argument for condition Lemma 3.1(iii) to be satisfied here gives the following remark.

**Remark.** *The issue of compactness or completeness of the bias spaces is rather a philosophical question. At best, one should think about the opportunity to "select" the subjective norm of reciprocity as a hypothetical choice rather than an alternative reflecting from a permanent conscious state of mind. Accordingly, the particular values of $\alpha, \beta$, and thus $\theta_1, \theta_2$, as emotional devices come into the conscious minds and initiate the actions only if the strategic situation requires it yet the whole spaces examining players' potentials are present at any time.*

In order to examine the viability of reciprocity, we will base our analyses on the results given with the Theorem of Appendix A and Lemma 3.1; however, we have to extend the setting somewhat since we assume the game $\Gamma_f$ to be asymmetric.

To emphasize the asymmetric character of the game consider now two different bias spaces with elements $\theta_1, \theta_2$ since $\gamma_1 \neq \gamma_2$, however symmetrical types are also possible, i.e. $\theta_1 = \theta_2$. So, $\theta_1 \in \Theta_1 = [-1 + \epsilon, 1 - \epsilon]$ and $\theta_2 \in \Theta_2 = [-1 + \epsilon, 1 - \epsilon]$ where $\theta_1 = \theta_2$ only if $\alpha \cdot \gamma_1 + (1 - \alpha) \cdot \gamma_2 = \beta \cdot \gamma_2 + (1 - \beta) \cdot \gamma_1$ with $\gamma_1 \neq \gamma_2$. Since the position of the players' roles is initially by no means (i.e. the players are either in position 1 or in position 2 with equal probability) it is relatively straightforward to construct an ex ante symmetric game setup. Thus, the profile parameter $\kappa = (\theta_1, \theta_2)$ is selected from the compact support $K = \Theta_1 \times \Theta_2 = [-1 + \epsilon, 1 - \epsilon] \times [-1 + \epsilon, 1 - \epsilon]$ and determines the evolving game parameter of the ex ante symmetric game with the distribution $G_t = (G_t^1, G_t^2)$ where $G_t^1$ corresponds to $\theta_1$ and $G_t^2$ to $\theta_2$; the $t$ will sometimes be dropped for convenience.[18] Then, the ex ante symmetric game payoff of an individual with type $\kappa = (\theta_1, \theta_2)$ competing with an individual of type $\tilde{\kappa} = \left( \tilde{\theta}_1, \tilde{\theta}_2 \right)$ is defined by

$$f(\kappa, \tilde{\kappa}) = \frac{f_1\left(\theta_1, \tilde{\theta}_2\right) + f_2\left(\theta_2, \tilde{\theta}_1\right)}{2}. \tag{3.8}$$

---

[18]Of course, the basic evolving trait is the norm of reciprocity, $\alpha$ (respective $\beta$), however, for the sake of clarity, it is sometimes benefiting to think of the overall concern as the evolving trait (consider $\theta_1, \theta_2$ as an initial random weighting of $\alpha$ and $\beta$).

Accordingly,

$$f(\kappa, G) = \frac{f_1(\theta_1, G^2) + f_2(\theta_2, G^1)}{2} \tag{3.9}$$

is the ex ante fitness to type $\kappa = (\theta_1, \theta_2)$ under the distribution $G = (G^1, G^2)$. Having this game texture allows us to transfer methodological results from Heifetz and Segev (2004) who use an asymmetric game setting which is close to ours in order to identify "the evolutionary role of toughness in bargaining" which gives the name to their essay. Extending the terminology of domination as in the Appendix A to the asymmetric game setting we have that $\tilde{\kappa} = \left(\tilde{\theta}_1, \theta_2\right)$ (or $\hat{\kappa} = \left(\theta_1, \widehat{\theta_2}\right)$) is dominated by $\kappa = (\theta_1, \theta_2)$ in iteration $n+1$ if for every $\acute{\kappa} = \left(\acute{\theta}_1, \acute{\theta}_2\right) \in U_n$ we have $f_1\left(\theta_1, \acute{\theta}_2\right) > f_1\left(\widetilde{\theta_1}, \acute{\theta}_2\right)$ (or $f_2\left(\theta_2, \acute{\theta}_1\right) > f_2\left(\widehat{\theta_2}, \acute{\theta}_1\right)$).

Accordingly, we reach:

**Lemma 3.2.** Dominance solvability of the asymmetric game $\Gamma_f$ with the two players' payoffs $f_1(\theta_1, \theta_2)$ and $f_2(\theta_2, \theta_1)$ implies dominance solvability of the ex ante symmetric game $\Gamma_{\Gamma_f}^\gamma$ with payoff $f(\kappa, \tilde{\kappa})$.

Using now the Theorem of Appendix A and Lemma 3.2, we have:

**Lemma 3.3.** If both players' asymmetric game is dominance solvable to $\theta_1^*(\alpha^*, \gamma_1, \gamma_2)$ and $\theta_2^*(\beta^*, \gamma_1, \gamma_2)$, respectively, then the profile $\kappa^* = (\theta_1^*, \theta_2^*)$ is the unique outcome of the ex ante symmetric game $\Gamma_{\Gamma_f}^\gamma$ under any regular and payoff-monotonic selection dynamics.

We are now able to prove our main result.

**Proposition 3.1.** Consider the game described above with $l \in [-1/4, 0) \cup (0, 3/5]$ and an extra requirement as given below the proof of this proposition (the requirement specifies the sets of $\gamma_1$ and $\gamma_2$ that we consider in different situations $l$). Then, any initial full-support of the distribution of biases $G = (G_0^1, G_0^2)$ will converge in distribution towards a unit mass on the pair of $(\theta_1^* = \alpha^* \gamma_1 + \beta^* \gamma_2, \theta_2^* = \beta^* \gamma_2 + \alpha^* \gamma_1)$

with the pair of reciprocity norms

$$\left( \alpha^* = 1/2 + 1/2 \, \frac{-4\,l - 4 + \sqrt{\Lambda\left(\gamma_1, \gamma_2, l\right)}}{l\left(l-2\right)\left(\gamma_1 - \gamma_2\right)}, \beta^* = 1/2 + 1/2 \, \frac{4\,l + 4 - \sqrt{\Lambda\left(\gamma_1, \gamma_2, l\right)}}{l\left(l-2\right)\left(\gamma_1 - \gamma_2\right)} \right),$$

with

$$\Lambda\left(\gamma_1, \gamma_2, l\right) = l^4 \left(\gamma_1 + \gamma_2 + 2\right)^2 - 4\,l^3 \left(\left(\gamma_1 + \gamma_2 + 1\right)^2 - 1\right) + $$
$$4\,l^2 \left(\left(\gamma_1 + \gamma_2 + 1\right)^2 - 1 - 4\,\gamma_1 - 4\,\gamma_2\right) + 16\,l\left(\gamma_1 + \gamma_2 + 2\right) + 16,$$

under any regular and payoff-monotonic dynamics.

*Proof.* According to the Theorem of Appendix A and Lemmata 3.1-3.3, the procedure of the proof is to find an equilibrium profile of the asymmetric game $\Gamma_f = (\{1, 2\}, \{\alpha, \beta\}, \{f_1, f_2\})$, and then check for sufficient conditions regarding dominance solvability. Thus, calculating first order conditions of $f_1\left(\theta_1, \theta_2\right)$ and $f_2\left(\theta_1, \theta_2\right)$ (cf. Eqs. (3.7)), i.e. $\frac{\partial f_1}{\partial \theta_1}\left(\theta_1, \theta_2\right) = 0$ and $\frac{\partial f_2}{\partial \theta_2}\left(\theta_1, \theta_2\right) = 0$, and solving for the biases yield

$$\theta_1 = -\frac{\left(\theta_2 l + l + 2\theta_2 + 2\right) l}{\theta_2 l^2 + l^2 - 2\theta_2 l - 2l - 4} \tag{3.10a}$$

$$\theta_2 = -\frac{\left(\theta_1 l + l + 2\theta_1 + 2\right) l}{\theta_1 l^2 + l^2 - 2\theta_1 l - 2l - 4}, \tag{3.10b}$$

where it is now reasonable to account for

$$\theta_1 = \alpha \cdot \gamma_1 + \left(1 - \alpha\right) \cdot \gamma_2 \tag{3.11a}$$

$$\theta_2 = \beta \cdot \gamma_2 + \left(1 - \beta\right) \cdot \gamma_1, \tag{3.11b}$$

in order to find equilibria on the level of reciprocity. To this end, we plug Eqs. (3.11) in Eqs. (3.10), equalize $\alpha$ and $\beta$, and solve for the equilibria.[19] Accordingly, we reach

$$\alpha^*_\pm = 1/2 + 1/2 \, \frac{-4\,l - 4 \pm \Phi\left(\gamma_1, \gamma_2, l\right)}{l\left(l-2\right)\left(\gamma_1 - \gamma_2\right)}$$

---

[19]Note that, mathematically, it does not matter whether we substitute the overall concern for its components in Eqs. (3.10) or in the fitness functions (Eqs. (3.7)).

and

$$\beta_\pm^* = 1/2 + 1/2\,\frac{4\,l + 4 \pm \Phi\left(\gamma_1,\gamma_2,l\right)}{l\left(l-2\right)\left(\gamma_1 - \gamma_2\right)},$$

where

$$\Phi\left(\gamma_1,\gamma_2,l\right) = \sqrt{\Lambda\left(\gamma_1,\gamma_2,l\right)},$$

with

$$\Lambda\left(\gamma_1,\gamma_2,l\right) = l^4\left(\gamma_1 + \gamma_2 + 2\right)^2 - 4\,l^3\left(\left(\gamma_1 + \gamma_2 + 1\right)^2 - 1\right) +$$
$$4\,l^2\left(\left(\gamma_1 + \gamma_2 + 1\right)^2 - 1 - 4\,\gamma_1 - 4\,\gamma_2\right) + 16\,l\left(\gamma_1 + \gamma_2 + 2\right) + 16,$$

where $\Lambda\left(\gamma_1,\gamma_2,l\right) > 0$ if $\gamma_1,\gamma_2 \in A = [-1+\epsilon, 1-\epsilon]$ and $l \in [-1/4, 1]$. Further, by regarding the restrictions on $\gamma_1,\gamma_2$ and the strategic setting $l$, and by analyzing the result sets of $\alpha_\pm^*$ and $\beta_\pm^*$, respectively, we find that $\alpha_-^*$ and $\beta_+^*$ are no possible solutions since $\alpha_-^* \notin [0,1]$ and $\beta_+^* \notin [0,1]$. In order to prove that $\alpha_-^*, \beta_+^* \notin [0,1]$ we have to consider 8 cases, or accordingly, we have to verify 8 conditions (each condition corresponds to one case), denote $(I)$ to $(VIII)$. The $\Phi\left(\gamma_1,\gamma_2,l\right)$ is dropped in the following since it is a positive value and of no account in any of the 8 conditions.

Case 1: Assume $\gamma_1 - \gamma_2 < 0$ and $l \in (0,1]$. Then, $\alpha_-^* < 0$ implies that $(I):= -4 < -l(l-2)(\gamma_1 - \gamma_2) + 4l$, which is true since the right hand side of $(I)$ is positive.

Case 2: Assume $\gamma_1 - \gamma_2 > 0$ and $l \in (0,1]$. Then, $\alpha_-^* > 1$ implies that $(II):= -4 < l(l-2)(\gamma_1 - \gamma_2) + 4l$, which is true since the right hand side of $(II)$ is positive.

Case 3: Assume $\gamma_1 - \gamma_2 < 0$ and $l \in [-0,25,0)$. Then, we analyze the same condition as under case 2, i.e. $(II){=}(III)$, but with negative $l$ and $\gamma_1 < \gamma_2$; however, this condition holds even under these constraints.

Case 4: Assume $\gamma_1 - \gamma_2 > 0$ and $l \in [-0,25,0)$. Then, we analyze the same condition as under case 1, i.e. $(I){=}(IV)$, but with negative $l$ and $\gamma_1 > \gamma_2$; however, this condition holds even under these constraints.

Case 5: Assume $\gamma_1 - \gamma_2 < 0$ and $l \in (0,1]$. Then, $\beta_+^* > 1$ implies that $(V):= 4 > l(l-2)(\gamma_1 - \gamma_2) - 4l$, which is true since $(V) = (-1) \cdot (I)$.

Case 6: Assume $\gamma_1 - \gamma_2 > 0$ and $l \in (0,1]$. Then, $\beta_+^* < 0$ implies that $(VI):=$

$4 > -l(l-2)(\gamma_1 - \gamma_2) - 4l$, which is true since $(VI) = (-1) \cdot (II)$.

Case 7: Assume $\gamma_1 - \gamma_2 < 0$ and $l \in [-0,25,0)$. Then, we analyze the same condition as under case 6 (respective 2), i.e. $(VII)=(VI)=(-1)\cdot(II)$, but with negative $l$ and $\gamma_1 < \gamma_2$; however, this condition holds even under these constraints.

Case 8: Assume $\gamma_1 - \gamma_2 > 0$ and $l \in [-0,25,0)$. Then, we analyze the same condition as under case 5 (respective 1), i.e. $(VIII)=(V)=(-1)\cdot(I)$, but with negative $l$ and $\gamma_1 > \gamma_2$; however, this condition holds even under these constraints.

Hence, the unique equilibrium profile to be further analyzed is given by

$$(\alpha_+^* = 1/2 + 1/2 \frac{-4\,l - 4 + \Phi(\gamma_1, \gamma_2, l)}{l\,(l-2)\,(\gamma_1 - \gamma_2)} = \alpha^*,$$
$$\beta_-^* = 1/2 + 1/2 \frac{4\,l + 4 - \Phi(\gamma_1, \gamma_2, l)}{l\,(l-2)\,(\gamma_1 - \gamma_2)} = \beta^*), \tag{3.12}$$

where we have dropped the subscripts "$+$" and "$-$" for convenience.

According to Lemma 3.1(i), the next step is to show that Eqs. (3.7) fulfill the properties of

($\blacktriangleleft$) "continuity",

($\blacktriangleright$) "twice differentiability", and

($\blacktriangledown$) "concavity",

with respect to $\alpha$ (respective $\beta$). By substituting $\theta_1$ and $\theta_2$ for its components, it is easy to comprehend that $\blacktriangleleft$ and $\blacktriangleright$ are satisfied. To show $\blacktriangledown$, review that we have defined the fixed and evolving dispositions by the restrictions

$$\gamma_1, \gamma_2 \in A = [-1 + \epsilon, 1 - \epsilon] \text{ with } \gamma_1 \neq \gamma_2,$$
$$\alpha, \beta \in B = [0, 1],$$

and the convex combinations, which totally identify the players, by

$$\left. \begin{array}{l} \theta_1 = \alpha \cdot \gamma_1 + (1 - \alpha) \cdot \gamma_2 \\ \theta_2 = \beta \cdot \gamma_2 + (1 - \beta) \cdot \gamma_1 \end{array} \right\} \in \Theta = [-1 + \epsilon, 1 - \epsilon]$$

Now, let us first check for concavity of Eqs. (3.7) in $\theta_1$ (respective $\theta_2$). Note that although we deal with asymmetric fitness functions, it is sufficient to show that one player's payoff is concave in the overall concern, i.e. $\frac{\partial^2 f_A}{(\partial \theta_A)^2}(\theta_A, \theta_B) < 0$ (for $A \in \{1, 2\}$ and $B = 3 - A$), since the pools of $\theta_A, \theta_B$ are equal. Calculating the second derivative with respect to Eqs. (3.7) yields

$$\frac{\partial^2 f_A}{(\partial \theta_A)^2}(\theta_A, \theta_B) = 2\, l^2 \frac{T(\theta_A)}{\left(l^2 + \theta_B l^2 - 4 + \theta_A l^2 + \theta_A l^2 \theta_B\right)^4},$$

where

$$
\begin{aligned}
T(\theta_A) = {}& -16 + \left(3\,\theta_B{}^2 + \theta_A + \theta_B{}^3 + 3\,\theta_B{}^2\theta_A + 3\,\theta_A\theta_B + 1 + \theta_A\theta_B{}^3 + 3\,\theta_B\right) l^5 \\
& + \left(16\,\theta_B + 14\,\theta_B{}^2 - 4\,\theta_B{}^2\theta_A - 2\,\theta_A\theta_B{}^3 + 6 - 2\,\theta_A\theta_B + 4\,\theta_B{}^3\right) l^4 \\
& + \left(-8\,\theta_B{}^2\theta_A - 8\,\theta_A - 16\,\theta_A\theta_B + 12\,\theta_B{}^2 + 24\,\theta_B + 12\right) l^3 \\
& + \left(-4\,\theta_B{}^2 - 8\,\theta_A\theta_B - 8\,\theta_A + 4\right) l^2 + (-16 - 16\,\theta_B) l
\end{aligned}
$$

Somewhat tedious calculations reveal that $T(\theta_A) < 0$ if $l \in [-1/4, 3/5]$ and $\theta_A, \theta_B \in \Theta$, and thus $\frac{\partial^2 f_A}{(\partial \theta_A)^2}(\theta_A, \theta_B) < 0$, if $l \in [-1/4, 0) \cup (0, 3/5]$ and $\theta_A, \theta_B \in \Theta$. Therefore, the players' best replies concerning their overall biases are given by Eqs. (3.10). However, whether the players' best replies concerning their reciprocity biases are implicitly given by Eqs. (3.10) is still an open question. Differentiating Eqs. (3.7) with respect to $\alpha$ (respective $\beta$) by applying the chain rule leads to

$$\frac{\partial^2 f_1}{(\partial \alpha)^2}(\theta_1, \theta_2) = \frac{\partial^2 f_1}{(\partial \theta_1)^2}(\theta_1, \theta_2)(\gamma_1 - \gamma_2)^2,$$

and

$$\frac{\partial^2 f_2}{(\partial \beta)^2}(\theta_1, \theta_2) = \frac{\partial^2 f_2}{(\partial \theta_2)^2}(\theta_1, \theta_2)(\gamma_2 - \gamma_1)^2,$$

where $(\gamma_1 - \gamma_2)^2 > 0$ and $(\gamma_2 - \gamma_1)^2 > 0$ in any case, and $\frac{\partial^2 f_1}{(\partial \theta_1)^2}(\theta_1, \theta_2) < 0$, $\frac{\partial^2 f_2}{(\partial \theta_2)^2}(\theta_1, \theta_2) < 0$ if $l \in [-1/4, 0) \cup (0, 3/5]$. Then $\frac{\partial^2 f_1}{(\partial \alpha)^2}(\theta_1, \theta_2) < 0$ and $\frac{\partial^2 f_2}{(\partial \beta)^2}(\theta_1, \theta_2) < 0$ under the same constraint concerning $l$.

The next step is to show that the slope of the best reply functions of the asymmetric bias game are less than 1 in absolute value (cf. Lemma 3.1(ii)). We derive the

two players' best reply functions concerning the reciprocity motive by calculating and solving the first order conditions of Eqs. (3.7) with respect to $\alpha$ (respective $\beta$), or equivalently, by plugging Eqs. (3.11) in Eqs. (3.10) and solving for the reciprocity norms. Accordingly, we reach

$$BR_1 = \alpha \left( \beta = \frac{\theta_2 - \gamma_1}{\gamma_2 - \gamma_1} \right) = \frac{\left(-l^2 + 2\gamma_2 l - 2l - \gamma_2 l^2\right)\theta_2 + 4\gamma_2 - \gamma_2 l^2 - l^2 + 2\gamma_2 l - 2l}{\left(l^2 - 2l\right)\left(\gamma_1 - \gamma_2\right)\theta_2 + \left(l^2 - 2l - 4\right)\left(\gamma_1 - \gamma_2\right)}$$

(3.13a)

$$BR_2 = \beta \left( \alpha = \frac{\theta_1 - \gamma_2}{\gamma_1 - \gamma_2} \right) = \frac{\left(\gamma_1 l^2 - 2\gamma_1 l + l^2 + 2l\right)\theta_1 + \gamma_1 l^2 - 2\gamma_1 l - 4\gamma_1 + l^2 + 2l}{\left(l^2 - 2l\right)\left(\gamma_1 - \gamma_2\right)\theta_1 + \left(l^2 - 2l - 4\right)\left(\gamma_1 - \gamma_2\right)}.$$

(3.13b)

The slopes of the best reply functions are given by

$$BR_1^s = \frac{d\alpha}{d\beta}(\beta) = -4\frac{l(2+l)}{\left(4 + (\gamma_1\beta - \gamma_1 - \beta\gamma_2 - 1)l^2 + (-2\gamma_1\beta + 2 + 2\gamma_1 + 2\beta\gamma_2)l\right)^2}$$

(3.14a)

$$BR_2^s = \frac{d\beta}{d\alpha}(\alpha) = -4\frac{l(2+l)}{\left(-4 + (\alpha\gamma_1 + 1 + \gamma_2 - \gamma_2\alpha)l^2 + (-2\alpha\gamma_1 - 2\gamma_2 - 2 + 2\gamma_2\alpha)l\right)^2},$$

(3.14b)

where

$$\sup_{l\in[-1/4,0)\cup(0,3/5]} \left| -\left(4l^2 + 8l\right) \right| \approx 6,24$$

and

$$\inf_{\substack{l\in[-1/4,0)\cup(0,3/5]\\ \gamma_1,\gamma_2\in A\\ \alpha,\beta\in B}} \left|\mathrm{den}\left(BR_1^s\right)\right| = \inf_{\substack{l\in[-1/4,0)\cup(0,3/5]\\ \gamma_1,\gamma_2\in A\\ \alpha,\beta\in B}} \left|\mathrm{den}\left(BR_2^s\right)\right| \approx 8,265,$$

with den $(\circ)$ symbolizing the denominator of $\circ$. Since $6,24 < 8,265$ the slopes of the best reply functions of the two players are less than 1 in absolute value, so condition Lemma 3.1(ii) holds.

In conclusion, by Lemmata (3.1-3.3), the ex ante bias game $\Gamma_{\Gamma_f}^\gamma$ is dominance solvable with $(\alpha^*, \beta^*)$ as in Eq. (3.12) as the unique Nash equilibrium profile, and the unique profile that survives any dynamic regular and payoff-monotonic process with full support on the one-dimensional reciprocity space. $\qquad\square$

In order to analyze only situations where the evolutionary outcome of reciprocity lies between 0 and 1, we need to compute the following requirement as announced

in Proposition 3.1.

**Requirement.** *The following relation is necessary to guarantee that $\alpha^*, \beta^* \in [0,1]$.*

$$|\Phi(\gamma_1, \gamma_2, l) - (4l + 4)| \leq |l(l - 2)(\gamma_1 - \gamma_2)|. \tag{3.15}$$

*Proof.* Since $\alpha^* + \beta^* = 1$, it is sufficient to show that $\alpha^* \in [0,1]$. Consider $0 \leq \alpha^* \leq 1$ with $\alpha^*$ as in Eq. (3.12), then we need to examine 2 cases:

Case 1: $l(l - 2)(\gamma_1 - \gamma_2) < 0$, i.e. $\text{sign}(l) = \text{sign}(\gamma_1 - \gamma_2)$, and

Case 2: $l(l - 2)(\gamma_1 - \gamma_2) > 0$, i.e. either $l < 0$ or $\gamma_1 < \gamma_2$.

Rearranging $0 \leq \alpha^* \leq 1$ given the first case leads to

$$4l + 4 + \underbrace{l(l - 2)(\gamma_1 - \gamma_2)}_{<0} \leq \Phi(\gamma_1, \gamma_2, l) \leq 4l + 4 - \underbrace{l(l - 2)(\gamma_1 - \gamma_2)}_{<0},$$

and the second case leads to

$$4l + 4 - \underbrace{l(l - 2)(\gamma_1 - \gamma_2)}_{>0} \leq \Phi(\gamma_1, \gamma_2, l) \leq 4l + 4 + \underbrace{l(l - 2)(\gamma_1 - \gamma_2)}_{>0}.$$

Subsuming both cases gives

$$4l + 4 - |l(l - 2)(\gamma_1 - \gamma_2)| \leq \Phi(\gamma_1, \gamma_2, l) \leq 4l + 4 + |l(l - 2)(\gamma_1 - \gamma_2)|,$$

and thus

$$|\Phi(\gamma_1, \gamma_2, l) - (4l + 4)| \leq |l(l - 2)(\gamma_1 - \gamma_2)|.$$

$\square$

Solving for a parameter of this expression does not give much additional insight, instead, in the Appendix B we show some representative situations that conform to this requirement (see Figure 3.3 and Table 3.1 there).

Having proved our basic result, it is relatively simple to derive a benchmark finding where the evolutionary players are not able to feel reciprocity.[20] Let us first define a one-dimensional population as follows.

---

[20]Qualitatively, the same result appears in an example of Heifetz et al. (2007a).

**Definition 3.2.** *A one-dimensional population here is a population as in the foregoing environment, but without reciprocity, and where the evolving trait is simply the intrinsic preference (the weight which is put on the opponents' material profit). That is, $\alpha = \beta = 1$ (cf. Definition 3.1), such that $\theta_A = 1 \cdot \gamma'_A + (1-1) \cdot \gamma'_B = \gamma'_A$ and $\theta_B = 1 \cdot \gamma'_B + (1-1) \cdot \gamma'_A = \gamma'_B$, where the "'" symbolizes "evolving" or "endogenized". Also, in this settig, we allow for $\gamma'_A = \gamma'_B$.*

Then, we reach the following result.

**Proposition 3.2.** Consider the one-dimensional population described above with a strategic interdependence according to $l \in [-1/4, 0) \cup (0, 3/5]$ and a mutation space given by $\theta_A = \gamma'_A, \theta_B = \gamma'_B \in \Theta = [-1 + \epsilon, 1 - \epsilon]$. Then, any initial full-support distribution of biases converges in distribution towards the unit mass on $l/(2-l)$, under any regular and payoff-monotonic dynamics.

*Proof.* As we have computed that $\frac{\partial^2 f_A}{(\partial \theta_A)^2}(\theta_A, \theta_B) < 0$ if $l \in [-1/4, 0) \cup (0, 3/5]$, it suffices to show that the slope of the best reply function is less than 1 in absolute value in this variant setting, since twice-differentiabilty and continuity, as also requested by Lemma 3.1, is obviously here. The best reply function of player $A$ is

$$BR_A = \theta_A(\theta_B) = \underset{\theta_A}{\operatorname{argmax}} f_A(\theta_A, \theta_B) = -\frac{(\theta_B l + l + 2\theta_B + 2)l}{\theta_B l^2 + l^2 - 2\theta_B l - 2l - 4}. \qquad (3.16)$$

The slope of the best reply function is

$$BR_A^s = \frac{d\theta_A}{d\theta_B}(\theta_B) = \frac{4(l+2)l}{(\theta_B l^2 + l^2 - 2\theta_B l - 2l - 4)^2}. \qquad (3.17)$$

Under assumptions $l \in [-1/4, 0) \cup (0, 3/5]$ and $\theta_A, \theta_B \in [-1 + \epsilon, 1 - \epsilon]$, we see that

$$\frac{dBR_A^s}{d\theta_B} = -8 \frac{(l+2)l(l^2 - 2l)}{(\theta_B l^2 + l^2 - 2\theta_B l - 2l - 4)^3} < 0.$$

Hence, Eq. (3.17) is decreasing in $\theta_B$, and thus maximized at $\theta_B = -1 + \epsilon$ and minimized at $\theta_B = 1 - \epsilon$. With $l \in [-1/4, 0)$, Eq. (3.17) is negative and the maximum absolute value occurs at $\theta_B = 1 - \epsilon$, where $|BR_A(\theta_B = 1 - \epsilon)| = \left| 4 \frac{(l+2)l}{(-2l^2 + l^2\epsilon + 4l - 2l\epsilon + 4)^2} \right| < 1$. With $l \in (0, 3/5]$, Eq. (3.17) is positive and the maximum

absolute value occurs at $\theta_B = -1 + \epsilon$, where $|BR_A(\theta_B = -1 + \epsilon)| = \left| 4 \frac{(l+2)l}{(l^2\epsilon - 2\,l\epsilon - 4)^2} \right| <$ 1. Since the variant game is dominance solvable, we find the outcome $l/(2 - l)$ by equalizing and solving for the biases with respect to Eq. (3.16), i.e. solving for $\theta_A$ in $\theta_A(\theta_B = \theta_A)$. $\qquad\square$

So, what is the significance of these results? By fielding this question, one should bear in mind that although the constraints of the equilibrium, which are determined by the game dynamical aspects and the model parameters, appear somewhat restricted, all assumption are intuitively plausible and of sufficient general character. The game dynamics allow for different initial populations to develop in the wide field of regularity and payoff monotonicity only provided that the population describes a compact interval in the line of reciprocity. Likewise, the payoff function which defines reproductive success in the society stands for a wide variety of different strategic games.

There are several observations which we can make easily. To start with, note that both players' equilibrium reciprocity values sum up to 1, i.e. $\alpha^* + \beta^* = 1$. This fact guarantees that both players' regard for the opponents' payoff is identical.

**Corollary 3.1.** $\theta_1^* = \theta_2^*$.

*Proof.* Since $\alpha^* + \beta^* = 1$, we have $\theta_1^* = \underbrace{\alpha^*}_{=1-\beta^*} \cdot \gamma_1 + \underbrace{\beta^*}_{=1-\alpha^*} \cdot \gamma_2 = \theta_2^*.$ $\qquad\square$

This result is not surprising because asymmetry of the equilibrium payoffs emerges not on the level of the players' overall concern but only with respect to the intrinisc preferences of the players. The next observation regards a comparing of our basic result with the outcome in the one-dimensional population model.

**Corollary 3.2.** *The two-dimensional population model may develop a different overall concern than the one-dimensional population model does.*

*Proof.* Let us write down only one simple example. Consider $l = -0,25$, $\gamma_1 = -0.8$, and $\gamma_2 = 0.4$. Then, the two-dimensional population develops an outcome that is

approximately $\alpha^* \cdot \gamma_1 + \beta^* \cdot \gamma_2 \approx -0,08$ and the outcome of the one-dimensional population model corresponds to $\frac{l}{2-l} = -0,\overline{1}$.  □

However, the fact that the strategic environment determines the sign of the overall concern is an observation which is of general character.

**Corollary 3.3.** *The players show a negative overall concern if strategic substitutes are present, i.e. $l < 0 \Rightarrow \theta_1^* = \theta_2^* < 0$. If the underlying game exhibits strategic complements, then, the players' value their opponents' payoff positively, i.e. $l > 0 \Rightarrow \theta_1^* = \theta_2^* > 0$.*

This result is reminiscent of the pioneering work of Bester and Güth (1998) where strategic complements leads to altruism and strategic substitutes to selfishness.[21] The following observation regards the reciprocity motive and its conclusion holds universally for the case of strategic substitutes, i.e. $l \in [-1/4, 0)$, and partly for the case of strategic complements ($l \in (0, 3/5]$).

**Corollary 3.4.** *There exists a threshold $\zeta(l) \in A$ such that*

$$\frac{\partial \alpha^*(\gamma_1, \gamma_2, l)}{\partial \gamma_2} < 0, \ for \ \gamma_1 < \zeta(l), \tag{3.18}$$

*and*

$$\frac{\partial \alpha^*(\gamma_1, \gamma_2, l)}{\partial \gamma_2} > 0, \ for \ \gamma_1 > \zeta(l), \tag{3.19}$$

*and for the second player likewise. A rough conclusion of this observation works as follows. If player A's intrinsic lies below the threshold, then the will to reciprocate to player B increases as player B gets more nasty. Likewise, if player A's intrinsic lies above the threshold, then the will to reciprocate to player B increases as player B gets more nice. Furthermore, the threshold is increasing in the strategic setting $l$.*

As the fractions of Ineq. (3.18) and Ineq. (3.19) are no continuous functions for the parameter constraints when strategies are complements, the conclusion derived

---

[21]However, the finding under strategic substitutes is restricted there, which is due to Bester and Güth's presumption of a non-negative preference space, and extends to spitefulness if the zero-barrier is abrogated (cf. Bolle, 2000; Possajennikov, 2000).

from these inequations is limited for $l \in (0, 3/5]$. In particular, for strategic comple-
ments and some values of $\gamma_1 \in A$, denote $\hat{\gamma}_1$, there is a critical point for $\gamma_2$ such that
$\alpha^*(\hat{\gamma}_1, \gamma_2 - \delta, l \in (0, 3/5]) < \alpha^*(\hat{\gamma}_1, \gamma_2, l \in (0, 3/5]) > \alpha^*(\hat{\gamma}_1, \gamma_2 + \delta, l \in (0, 3/5])$, for
positive $\delta$, and for the second player likewise.

## 3.4    Conclusion

Building upon the assumption that individuals adjust their actions to achieve higher
subjective utility, while dynamical pressures change the composition of reciprocal
preferences in the population according to the players' objective gains, this study
provides a prognose concerning the emergence of strong reciprocity in a wide class of
strategic interaction. The specific conception of reciprocity is defined by the players
tendency to response to the perceived intrinsic attitudes of others. The basic finding
is that a high degree of flexibility (in the reciprocal sense) pays off. In our setting,
it is the strategic environment and the specific other players' type which determine
the players' behavior and only marginally their usual, exogeneous given, intrinsic
attitudes. The unique dominance solvable profile of reciprocity in our asymmet-
ric setting, and hence the only survivor under any regular and payoff monotonic
selection process, motivates an altruistically inclined player to behave spitefully if
strategies are substitutes and the opponent player is intrinsically spiteful. Con-
versely, the reciprocity norm of a spiteful player motivates him to show altruistic
behavior if strategies are complements and the opponent player is intrinsically al-
truistic. In this regard, our study provides an economic and cultural explanation to
the question why people often show a different behavior than they usually would.
The study further substantiates related work which shows that the strategic envi-
ronment determines the players equilibrium behavior in one-dimensional preference
models. In particular, the fact that strategic substitutes lead to a negative emotion
and strategic complements to a positive emotion is once more confirmed, even in

the new setting of two preference dimensions.

A further conclusion of our work is that if player A is relatively spiteful, then player A's tendency to reciprocity is higher the nastier player B—if player A is relatively nice, then player A's tendency to reciprocity is higher the nicer player B. However, this conclusion is restricted in the sense that it holds universally only for the case of strategic substitutes.

One can think of several extensions of our basic model. Since preference biases act like commitment devices which form the other players' equilibrium strategies to a certain extent, it would be interesting to explore whether the qualitative results maintain in cases where the players do not recognize the other players' types perfectly. For instance, one can think of situations where the players' types are observed with some noise because they do not learn the types of each other, or where some players intentionally signal a wrong disposition in order to benefit from the resulting strategic effect. It would further be interesting whether our results can be identified in experimental studies—an admittedly subtle task since the individual usual types have to be ascertained before a norm of reciprocity can be examined.

## 3.5 Appendices

### 3.5.1 Appendix A. A Generic Class of Evolutionary Dynamics

To analyze the evolutionary viability of reciprocity, we build upon the generic class of selection dynamics as proposed by Heifetz et al. (2007a). This section is intended to sketch the distinctive attributes and advantages of this class. To this end, imagine the following.

At each instant in time $t \geq 0$, two players are randomly drawn from a continuum population to play a certain game. In particular, the population is characterized by the distribution $G_t \in \Delta(\Theta)$, where $\Delta(\Theta)$ is the set of Borel probability distributions over the compact space $\Theta = \left[\underline{\theta}, \bar{\theta}\right]$. The population evolves over time in the space of $\Delta(\Theta)$ according to the following differential equation.

$$\dot{G}_t(S) = \int_S g(\theta, G_t) \, dG_t(\theta), \ S \subseteq \Theta \text{ Borel measurable,} \qquad \text{(A.3.1)}$$

where $g : \Theta \times \Delta(\Theta) \to \mathbb{R}$ is a continuous growth-rate function. The following definition further specifies the dynamics.

**Definition 3.3.** *The continuous growth-rate function $g : \Theta \times \Delta(\Theta) \to \mathbb{R}$ is payoff-monotonic and regular if for every $G \in \Delta(\Theta)$, the following conditions hold:*

- *A higher average fitness corresponds to a higher growth-rate, or formally,*

$$\int f\left(\theta, \acute{\theta}\right) dG_t\left(\acute{\theta}\right) > \int f\left(\tilde{\theta}, \acute{\theta}\right) dG_t\left(\acute{\theta}\right) \iff g(\theta, G_t) > g\left(\tilde{\theta}, G_t\right). \qquad \text{(A.3.2)}$$

- *$G_t$ is a probability distribution for every $t$,*

$$\int_\Theta g(\theta, G) \, dG(\theta) = 0. \qquad \text{(A.3.3)}$$

- *g can be extended to the domain $\Theta \times X$, where $X$ is the set of signed Borel measures with variational norm smaller than 2, such that $g$ is uniformly bounded and Lipschitz continuous on $\Theta \times X$. Formally,*

$$\sup_{\theta \in \Theta} \left| g\left(\theta, G_t\right) \right| < \infty$$

$$\sup_{\theta \in \Theta} \left| g\left(\theta, G_t\right) - g\left(\theta, \widetilde{G}_t\right) \right| < K \left\| G_t - \widetilde{G}_t \right\|, G_t, \widetilde{G}_t \in X, \tag{A.3.4}$$

*for some constant $K$, where $\|G\| = \sup_{|h| \leq 1} \left| \int_\Theta h \, dG \right|$ is the variational norm on signed measures.*

Oechssler and Riedel (2001, Lemma 3) show that regularity of $g$ guarantees that the mapping $G \to \int_\Theta g\left(\cdot, G\right) dG$ is bounded and Lipschitz continuous in the variational norm, which implies that the differential Eq. (A.3.1) has a unique solution for any initial distribution $G_0$.[22]

The dynamics defined here formalize the simple idea that only individuals who play well in the population increase while individuals who play badly decrease. As mentioned in the introduction, the underlying evolving process may rely on a biological level or on a cultural level. Accordingly, more successful types have more descendantes who carry the genes of their parents, or more successful types are more likely to be adopted under a cultural process of education or imitation from role-models. Alternatively, Heifetz et al. (2007a) mention that the same mathematical

---

[22]To see that the replicator dynamics forms a special case of the generic dynamics determined by Eq. (A.3.1), consider a growth-rate function which evolves according to the subtraction of the population average success from the success of a single type (the difference is sometimes called the *excess payoff*). Formally, $G_t$ evolves according to

$$\dot{G}_t\left(S\right) = \int_S \left[f\left(\theta_i, G_t\right) - f\left(G_t, G_t\right)\right] dG_t\left(\theta_i\right), \ S \subseteq \Theta,$$

where

$$f\left(\theta_i, G_t\right) = \int_\Theta f_i\left(\theta_i, \theta_j\right) dG_t\left(\theta_j\right)$$

is the expected success of individual $i$ played against a randomly chosen individual $j$. And, the population average success is

$$f\left(G_t, G_t\right) = \int_\Theta \int_\Theta f_i\left(\theta_i, \theta_j\right) dG_t\left(\theta_j\right) dG_t\left(\theta_i\right).$$

structure is compatible with the idea that successful individuals have more influence on the dynamic process since they appear more often in economic interactions, and so, are more likely to be reproduced.

Dealing with the concept of dominance solvability requires a deeper understanding of the concept of domination and some additional notations. To begin with, we say that $\theta'$ is dominated by $\theta$ whenever $f(\theta, \theta'') > f(\theta', \theta'')$ for every $\theta'' \in \Theta$. Then, let $D_1$ denote the set of types $\theta'$ which are dominated by some $\theta \in \Theta$, and $U_1$ is the set of undominated types, i.e. $U_1 = \Theta \setminus D_1$. Further, $D_n$ is the set of dominated types after at most $n$ iterations and the set of undominated types is accordingly $U_n = \Theta \setminus D_n$. Then, $\theta' \in U_n$ is dominated in iteration $n+1$ by $\theta \in U_n$ if $f(\theta, \theta'') > f(\theta', \theta'')$ for every $\theta'' \in U_n$. We say that $\theta'$ is *serially dominated* if it is dominated after some number of iterations. Under regular and payoff monotonic dynamics any serially dominated types are extinct in the limiting population. A game is dominance solvable if there is a unique type that is not serially dominated. The analyzes of reciprocity in this paper is based on this class of selection dynamics and on the following theorem of Heifetz et al. (2007a) which extends results from Samuelson and Zhang (1992), where the population evolves according to matrix games, to continuous strategy spaces.

**Theorem.** *(Heifetz et al. (2007a)) Consider a symmetric two-player game with strategy space $\Theta = [\underline{\theta}, \overline{\theta}] \subset \mathbb{R}$, a continuous payoff function $f : \Theta \times \Theta \to \mathbb{R}$, and a regular, payoff-monotonic growth-rate function $g : \Theta \times \Delta(\Theta) \to \mathbb{R}$. Moreover, assume that the population $G$ has initially full support over the compact space $\Theta$ and evolves according to the differential equation as defined by Eq. (A.3.1). Then, types $\theta$ which are serially dominated are asymptotically weeded out, i.e. they have a neighborhood $V \ni \theta$ for which $\lim_{t \to \infty} G_t(V) = 0$. In particular, when the game is dominance solvable to equilibrium $\theta^*$, then $G_t$ converges in distribution to a unit mass at $\theta^*$.*

## 3.5.2 Appendix B. Figure and Table

Figure 3.3 presents the possible sets of $\gamma_1$ and $\gamma_2$ in 5 different strategic situations $l$: $l = -0,25$, $l = -0,05$, $l = 0,05$, $l = 0,25$, and $l = 0,5$. The graphics show that we analyze situations where the intrinsic preferences of the two populations are predominantly different, i.e. $\text{sign}(\gamma_1)$ is predominantly different from $\text{sign}(\gamma_2)$. The graphics also show that in the case of strategic substitutes $l < 0$, we do not consider populations with positive $\gamma_1$ and positive $\gamma_2$; likewise, in the case of strategic complements, we do not consider cases where $\gamma_1$ and $\gamma_2$ are both negative.

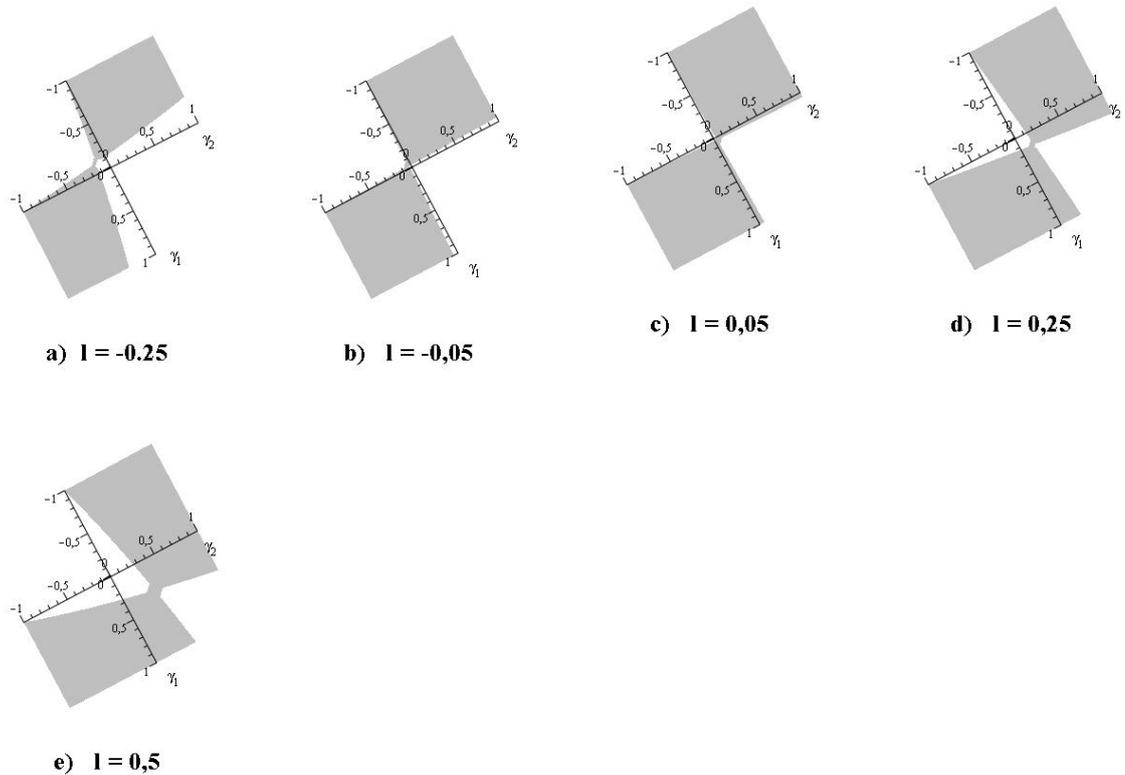a) l = -0.25  b) l = -0,05  c) l = 0,05  d) l = 0,25

e) l = 0,5

Figure 3.3: $\gamma_1, \gamma_2$-sets.

Table 3.1 shows some dominance solvable results with respect to different game

parameters. For expositional clarity we focus on the same "representative" environments as in Figure 3.3: we assume 5 different cases which examine the strategic environment, i.e. $l = -0,25$, $l = -0,05$, $l = 0,05$, $l = 0,25$, and $l = 0,5$, each with a two population setting of different intrinsic preferences. The table basically shows that a dominance solvable reciprocity trait is apt to reverse the players intrinsic attitude. This finding holds in settings with strategic substitutes ($l < 0$) as well as in those with strategic complements ($l > 0$). For example take $l = -0,25$, i.e. strategic substitutes, then if the $\gamma_1$-player is moderate altruistic ($\gamma_1 = 0,3$) and the opponent player is of type $\gamma_2 = -0,3$ then the dynamics drive the $\gamma_1$-player to a relatively reciprocal norm $(0,25)$ such that the overall concern becomes negative. The fact that the sign of the strategic environment determines the sign of the overall concern can also be observed in the table.

| Exogeneous game parameters | | | Dominance solvable traits (approx. values) | | |
|---|---|---|---|---|---|
| Strategic setting | Intrinsic preferences | | Reciprocity | | Overall concern |
| $l$ | $\gamma_1$ | $\gamma_2$ | $\alpha^*$ | $\beta^*$ | $\theta_1^* = \theta_2^*$ |
| $-0,25$ | $-0.9$ | $0,1$ | $0,125$ | $0,875$ | $-0.025$ |
| $-0,25$ | $0,3$ | $-0,3$ | $0,25$ | $0,75$ | $-0,15$ |
| $-0,05$ | $-0,6$ | $0,3$ | $0,35$ | $0,65$ | $-0,015$ |
| $-0,05$ | $0,6$ | $-0,1$ | $0,087$ | $0,913$ | $-0,039$ |
| $0,05$ | $0,7$ | $-0,2$ | $0,262$ | $0,738$ | $0,036$ |
| $0,05$ | $-0,7$ | $0,5$ | $0,4$ | $0,6$ | $0,02$ |
| $0,25$ | $0,4$ | $-0,5$ | $0,67$ | $0,33$ | $0,103$ |
| $0,25$ | $-0,3$ | $0,2$ | $0,19$ | $0,81$ | $0,105$ |
| $0,5$ | $0,6$ | $0,1$ | $0,478$ | $0,522$ | $0,339$ |
| $0,5$ | $-0,7$ | $0,3$ | $0,166$ | $0,834$ | $0,134$ |

Table 3.1: Results in different situations.

# Chapter 4

# Overconfidence in Tullock Contests: An Evolutionary Approach

*We explore evolutionarily stable levels of self-confidence in Tullock contests for finite and infinite populations. While the players match exactly their true value of self-confidence if the population is infinite, they always exhibit a tendency to overconfidence if the population is of finite size. More precisely, the smaller the size of the population, the stronger the evolutionarily stable degree of overconfidence. Methodologically, we use the approach of indirect evolution where players maximize given preferences which evolve according the evolutionary fitness they induce. We further establish a conformity between evolutionarily stable effort as calculated under direct evolution, and the equilibrium effort which is induced by the evolutionarily stable level of overconfidence given indirect evolution.*

## 4.1 Introduction

The deviations from standard routines of the homo economicus model have recently revolutionized the branch of behavioral economics. Economic decision makers are not as selfish, rational, and perfectly informed as presupposed in much traditional economic work. Due to observations of real life interaction and experimental economics, the consideration of "imperfect" economic agents occurs by now in many different areas of economics. In many experiments, people show "irrational" and/or "unselfish" behavior, for instance by "burning money" in order to either punish or reward others or for other reasons.[1]

In the present work, we consider contests where players make irrecoverable investments to influence their probability of winning. Evidently, this sort of competition is a very common phenomenon, and the related literature include, among others, litigation (Robson and Skaperdas, 2008), R&D competition (Nalebuff and Stiglitz, 1983), sporting competition (Szymanski, 2003), lobbying (Baye et al., 1993) and rent-seeking (Tullock, 1980). In particular, we investigate the probably most studied form of contests, namely Tullock contest, and achieve to substantiate detected deviations which occur from predictions of rational Nash equilibrium theory on the one side and results from experimental economics on the other. The brainchild of Tullock (1980) is that the individual probability of winning a (howsoever designed) contest is an increasing function of the own effort. In Nash equilibrium theory, it is established that the accumulated efforts of the contest participants never exceed the rent. However, experimental contests show that people may overexpenditure in relation to risk-neutral Nash expenditure levels.[2] ESS theory and direct evolution of effort levels already gives an answer to this puzzle since overdissipation can be an ESS outcome of a Tullock contest as shown by Hehenkamp et al. (2004).

The main objective of the present study is to further identify evolutionary ratio-

---

[1]Cf. Camerer (2003) for a great overview of the experimental literature.
[2]Cf. Hörisch and Kirchkamp (2010) and Morgan et al. (2008), and the references therein.

nales for departures from Nash equilibrium effort levels by allowing the players to exhibit an idiosyncratic feeling of self-confidence, i.e. the players are able to perceive underconfidence and overconfidence.[3] The idea to integrate the dimension of self-confidence in the Tullock contest environment builds upon recent empirical findings that overconfident subjects are more likely to self-select into more competitive settings than unbiased individuals (cf. Dohmen and Falk, 2006). Methodologically, we use indirect evolution where subjective and objective payoffs are two different values; the first is pursued by the players who exhibit an idiosyncratic parameter of self-confidence and the second measures the evolutionary success of this idiosyncratic parameter (cf. Güth et al., 2002). Indirect evolution is by now the central tool in behavioral economics and, in particular, evolutionary game theory, to explain distortions from the homo economicus model in experiments and real life. The methodology has recently stirred the contest literature. For example, Konrad (2004) models an all-pay auction (a limiting case of the Tullock contest where the player with the highest expenditure wins for certain) with two types of individuals: altruistic and envious ones. He considers incomplete information and an infinitely repeated game, and identifies an "interior equilibrium", i.e. both types exist in equilibrium such that they form a symbiosis. Schmidt (2009) considers a related model but with full information and he uses a standard Tullock contest. He finds that a population of altruists performs better than a population of envious players but an envious player has the edge over an altruist if both compete against each other. Moreover, he finds an evolutionary advantage of the envious player in that very altruistic players always die out but very envious players do so only under certain conditions. The result that "negatively interdependent" players (in the sense of suffering if others win) feature an evolutionary advantage can also be found in Leininger (2009). Methodically, he uses indirect evolution and Schaffer (1988)'s ESS version for finite populations to prove that evolutionarily stable prefer-

---

[3]If we use the expression "self-confidence", we mean its dimension which ranges from strong underconfidence over the neutral ("true") level to strong overconfidence.

ences exhibit negative interdependence which rationalizes higher expenditures than standard Nash behavior without distortions. Precisely, he finds that evolutionarily stable spite under indirect evolution yields to the same more aggressive behavior as calculated with the ESS conception under direct evolution. The build-up of our approach is very close to Leininger's—however, instead of altruism/spite we explore the level of self-confidence. As Leininger identifies spite as an evolutionary rationale for aggressive behavior in contests, we find overconfidence as another evolutionary rationale. Interestingly, the equivalence result of Leininger applies to the same extent in our setting, too. That is, direct evolution of effort and indirect evolution of self-confidence guides to the same behavior.[4]

The remainder of this paper is organized as follows. In the next section, we review a standard Tullock contest and corresponding Nash and ESS behavior, respectively. In section 4.3, we introduce the dimension of self-confidence in the Tullock environment, and detect evolutionarily stable degrees for finite and infinite populations. We further establish an equivalence result between behavior as induced by evolutionarily stable self-confidence given indirect evolution and the effort level which is calculated under direct evolution. We conclude with section 4.4, summarizing what has been learned.

## 4.2 Nash and ESS Behavior in Standard Tullock Contests

In order to prepare for the analysis of the next section, we first review a standard Tullock contest as well as corresponding Nash equilibrium and ESS behavior, respectively. The description of Nash and ESS behavior draws on Leininger (2009) and Hehenkamp et al. (2004). We treat the results as a benchmark to those which occur in the subsequent research.

---

[4]In the same way, Boudreau and Shunda (2010) analyze price perceptions in Tullock contests and identify the same equivalence result.

### 4.2.1 Tullock Contest and Nash Behavior

A standard Tullock contest works as follows. Consider the (expected) payoff of individual $i$, $i = 1, ..., n$:

$$\Pi_i = \mathcal{P}_i V - c_i(e_i), \tag{4.1}$$

where we resort to the following often used assumptions:

$$\mathcal{P}_i = \begin{cases} \frac{e_i^r}{\sum_{j=1}^{n} e_j^r} & \text{if } \max\{e_1, ..., e_n\} > 0 \\ 1/n & \text{otherwise} \end{cases}, \tag{4.2}$$

i.e., the so-called "contest success function" $\mathcal{P}_i$ is of the usual logit form, and we further assume currently that

$$c_i(e_i) = e_i, \tag{4.3}$$

meaning player $i$'s costs $c_i(e_i)$ and effort $e_i \in \mathbb{R}_0^+$ are equivalent. In the following, we refer to the corresponding payoff function

$$\pi_i = \mathcal{P}_i V - e_i, \tag{4.4}$$

as the *objective* (or material) payoff function (in the sense of intersubjective measurable). Recall that $r \in \mathbb{R}_0^+$ in the contest success function is called the technology parameter which represents the implemented efficiency of the contest; it is also termed the discriminatory power of the contest and determines the influence of a players' effort on his probability of winnning.[5] As is well-known, the unique Nash equilibrium effort, where all players maximize their material profits, correspond to

$$e^* = \frac{n-1}{n^2} r V, \tag{4.5}$$

given that $r \leq n/(n-1)$.

An important element of the Tullock contest is the dissipation rate. The dissipation rate determines the part of the price which is spent by the players, and is

---

[5]For example, $r = 1$ is a lottery, $r = 0$ eliminates the impact of the players' efforts on their winning probabilities, and $r \to \infty$ is an all-pay auction where the player with the highest effort wins for certain.

formally defined via the ratio

$$\mathcal{D} = \frac{\sum_{j=1}^{n} e_j}{V}. \tag{4.6}$$

In equilibrium it holds that $\sum_{j=1}^{n} e_j^* = ne^* = \frac{n-1}{n}rV$, such that $\mathcal{D}^* = \frac{n-1}{n}r$. As overdissipation (full or underdissipation) occurs if $\mathcal{D} > (= \text{or} <)$ 1, respectively, one can conclude that overdissipation is incompatible with rational Nash behavior since $\mathcal{D}^* = \frac{n-1}{n}r \leq 1$ given that $r \leq n/(n-1)$. The result that overdissipation do not occur in equilibrium with rational agents is robust to many modifications of the Tullock environment.[6] A big puzzle in this regard is that Nash behavior evidently fails to be a proper prediction since experimental results often show overdissipation of the players (cf. Morgan et al., 2008; Hörisch and Kirchkamp, 2010, and each with the references therein).

### 4.2.2 ESS Behavior

As is well-known an ESS, for evolutionarily stable strategy, is a robust strategy such that if all members of a group adopt it there is no other minority-strategy that performs better in expectation (cf. Maynard-Smith and Price, 1973). The following notion is standard for games with infinite populations.

**Definition 4.1.** $e^{ESS}$ *is evolutionarily stable iff, for all $e \neq e^{ESS}$, the equilibrium condition*

$$(i) \qquad \pi(e^{ESS}, e^{ESS}) > \pi(e, e^{ESS})$$

*or, in the case of equity, the stability condition*

$$(ii) \qquad \pi(e^{ESS}, e^{ESS}) = \pi(e, e^{ESS}), \pi(e, e) < \pi(e^{ESS}, e)$$

*holds.*

---

[6]For example, to models with loss aversion (Cornes and Hartley, 2003) or risk aversion (Konrad and Schlesinger, 1997). Cf. the references in Baharad and Nitzan (2008) for further examples.

Hehenkamp et al. (2004) use Schaffer (1988)'s ESS version for finite populations (where ESS is *not necessarily* a refinement of Nash's concept), in the Tullock contest environment and, interestingly, find that ESS behavior yields to more aggressive play than Nash behavior. The methodology is all important for our analysis. To review this result and to prepare for our analysis, consider a population of finite size $N$ where only $n \subseteq N$ players participate in the contest. Those $n$-participants are chosen randomly and with equal probability. Then, consider the following definition.

**Definition 4.2.** *Let a contest-strategy $e$ be adopted by all players $i$, $i = 1, ..., n$. A mutant strategy $\bar{e} \neq e$ can invade $e$, if the payoff of a single player with strategy $\bar{e}$ (against $e$ of the $(n-1)$ other players) is strictly higher than the payoff of a player with strategy $e$ (against the single mutant with strategy $\bar{e}$ and $(n-2)$ other players with strategy $e$). A strategy $e^{ESS}$ is evolutionarily stable, if it cannot be invaded by any mutant strategy.*

Assume now in the Tullock contest environment that player 1 is the only mutant player with strategy $\bar{e}$ such that all $(n-1)$ other players use strategy $e$. Then, the expected payoff of the mutant player is

$$\pi_1(\bar{e}, e, ..., e) = \frac{\bar{e}^r}{(n-1)e^r + \bar{e}^r}V - \bar{e}, \tag{4.7}$$

and the expected payoff of one of the $n-1$ players $i$, $i \in \{2, ..., n\}$, with effort $e$ is

$$\pi_i = \left(1 - \frac{n-1}{N-1}\right)\pi_i(e, ..., e) + \frac{n-1}{N-1}\pi_i(\bar{e}, e, ..., e), \tag{4.8}$$

since the probability that an ESS player face the mutant is $(n-1)/(N-1)$. In order to find an ESS, we reach the following maximization problem:

$$\max_e \pi_1(e, e^{ESS}, ..., e^{ESS}) - (1 - \frac{n-1}{N-1})\pi_i(e^{ESS}, ..., e^{ESS}) - \frac{n-1}{N-1}\pi_i(e, e^{ESS}, ..., e^{ESS}),$$

and thus, by dropping the constant second term:

$$\max_e \pi_1(e, e^{ESS}, ..., e^{ESS}) - \frac{n-1}{N-1}\pi_i(e, e^{ESS}, ..., e^{ESS}). \tag{4.9}$$

Accordingly, we recognize that the maximization problem is not only about pursuing higher own payoff but also about lowering others' payoff. One can refer to this type of behavior as "spite" (cf. Hamilton, 1971). Provided that $r \leq n/(n-1)$, the unique ESS-effort accords to

$$e^{ESS} = \frac{(n-1)N}{(N-1)n^2}rV,$$

(cf. Hehenkamp et al., 2004, Theorem 5). As $\lim_{N \to \infty} e^{ESS} = \lim_{N \to \infty} \frac{(n-1)N}{(N-1)n^2}rV = e^* = \frac{n-1}{n^2}rV$, we conclude that the difference of Nash and ESS behavior dissapears if the population is infinite. Further, given that $N = n$, it holds that ESS behavior is more aggressive than Nash behavior since $\frac{r}{n}V > \frac{n-1}{n^2}rV$. Unlike under Nash induced efforts, aggregate ESS effort may over-dissipate the rent and, furthermore, does not depend on the number of players but only on the given contest technology and the given rent. In particular, $\mathcal{D}^{ESS} = n \cdot e^{ESS} = r \cdot V$, and thus over-dissipation occurs if $r > 1$, under-dissipation occurs if $r < 1$, and full-dissipation occurs if $r = 1$.

## 4.3    Overconfidence: The Distorted Case

Following the indirect evolutionary approach of Güth and Yaari (1992), we differentiate between subjective payoff (or utility), which determines the players' idiosyncratic biases on the dimension of self-confidence, and objective payoff, which measures the evolutionary fitness of the particular bias. As noted in the foregoing section, the objective payoff $\pi_i$ is given by the standard expected Tullock success with costs $c_i(e_i) = e_i$, cf. Eq. (4.4). The subjective payoff is described in the following subsection.

### 4.3.1    Subjective Utility

We assume in the following that the contest participants may differ in the way they perceive their ability in the Tullock contest. More precisely, we assume that the

players have different perceptions about their effort costs (Ludwig et al., 2010, use a similar model to analyze overconfidence in Tullock contests; however, not from an evolutionary point of view):

$$c_i(e_i) = (1 + \chi_i) e_i, \tag{4.10}$$

where $\chi_i \in \mathcal{T}_i = \mathbb{R}$ identifies player $i$'s idiosyncrasy on the dimension of self-confidence. We interpret a person with type $\chi_i < (>) 0$ as overconfident (under-confident) since individual $i$ perceives his effort cost lower (higher) than it really is, $(1 + \chi_i) e_i < (>) e_i$ with $\chi_i < (>) 0$. Accordingly, the subjective utility of player $i$ is given by:

$$\mathcal{U}_i = \mathcal{P}_i V - c_i(e_i) = \pi_i - \chi_i e_i = \mathcal{P}_i V - (1 + \chi_i) e_i. \tag{4.11}$$

Hence, a player's subjective utility ($\mathcal{U}_i$) differs from his objective payoff ($\pi_i$) to the extent of $|\chi_i e_i|$. The players strive to maximize their subjective utility with an adequate effort level and we denote the corresponding effort game by $\Gamma = (I, \mathcal{S}_i, \mathcal{U}_i)$, where $I$ is the set of contest participants (players) with biases $\chi_i$ that are randomly drawn from the population $N$ to play the contest game, and $\mathcal{S}_i = \mathbb{R}^+ \ni e_i$. The assumptions made in this subsection are binding for the following subsections.

### 4.3.2 A Preliminary Result

The following example is shown to rationalize overconfidence in a Tullock contest environment without going in evolutionary details. Consider a two-player Tullock contest, i.e. $I = \{1, 2\}$, and a technology parameter as in a lottery, i.e. $r = 1$. In this game, both players maximize their payoff function but they differ in the way they perceive their payoffs; the first player is a profit maximizer while the second has a feeling of self-confidence. Formally,[7]

$$\acute{\mathcal{U}}_1 = \pi_1 = \frac{e_1}{e_1 + e_2} V - e_1,$$

---

[7]The ´ symbolizes the belonging to our motivating example of this subsection.

and

$$\acute{\mathcal{U}}_2 = \pi_2 - \chi_2 e_2 = \frac{e_2}{e_1 + e_2} V - (1 + \chi_2) e_2.$$

The Nash equilibrium efforts are then given by the profile

$$\left(\acute{e}_1^*, \acute{e}_2^*\right) = \left(\frac{\chi_2 V}{(1 + \chi_2)^2}, \frac{V}{(1 + \chi_2)^2}\right),$$

which yield subjective payoffs (recall that the subjective payoff of player 1 is identical with his objective payoff)

$$\acute{\mathcal{U}}_1^* = \frac{V \chi_2^2}{(1 + \chi_2)^2},$$

and

$$\acute{\mathcal{U}}_2^* = \frac{V}{(1 + \chi_2)^2}.$$

Hence, a slight disposition of player 2 is subjective beneficial in relation to the subjective utility of player 1 since $\acute{\mathcal{U}}_2^* > \acute{\mathcal{U}}_1^*$ if $\chi_2 \in (-1, 1)$. If the disposition of player 2 is too strong, then player 1 is subjective better of with his disability to perceive a distorted degree of self-confidence. However, due to the indirect evolutionary game which we suppose, only the objective payoff is decisive for determinig players' success. The objective payoffs are

$$\acute{\pi}_1^* = \frac{V \chi_2^2}{(1 + \chi_2)^2},$$

and

$$\acute{\pi}_2^* = \frac{\chi_2 V}{(1 + \chi_2)^2}.$$

Note that player 2 earns exactly the equilibrium effort of player 1. Provided that $\chi_2 \in (0, 1)$, the "distorted" player 2 is also the more successful one in evolutionary terms. However, the puzzle regarding which type is more successful concerning evolutionary stability is not solved yet. We address this issue in the following subsection by allowing both players to perceive the dimension of self-confidence.

### 4.3.3 Main Results

Consider again the case $I = \{1, 2\}$ and $r = 1$. However, we suppose now that both players are able to recognize the dimension of self-confidence. The subjective payoffs are

$$\mathcal{U}_1 = \mathcal{P}_1 V - c_1\left(e_1\right) = \pi_1 - \chi_1 e_1 = \frac{e_1}{e_1 + e_2} V - \left(1 + \chi_1\right) e_1, \qquad (4.12)$$

and

$$\mathcal{U}_2 = \mathcal{P}_2 V - c_2\left(e_2\right) = \pi_2 - \chi_2 e_2 = \frac{e_2}{e_1 + e_2} V - \left(1 + \chi_2\right) e_2. \qquad (4.13)$$

As is the standard procedure under the indirect evolutionary approach, both players maximize their subjective payoff which may be due to some learning process. It can easily be calculated that the unique Nash equilibrium profile of the effort game $\Gamma = \left(\{1, 2\}, \mathcal{S}_{i=1,2}, \mathcal{U}_{i=1,2}\right)$, where the players maximize $\mathcal{U}_1$ and $\mathcal{U}_2$, respectively, is then given by

$$\left(e_1^*, e_2^*\right) = \left(\frac{v\left(1 + \chi_2\right)}{\left(\chi_1 + \chi_2 + 2\right)^2}, \frac{v\left(1 + \chi_1\right)}{\left(\chi_1 + \chi_2 + 2\right)^2}\right). \qquad (4.14)$$

The subjective equilibrium payoffs are

$$\mathcal{U}_1\left(e_1^*, e_2^*\right) = \mathcal{U}_1^* = \frac{v\left(1 + \chi_2\right)^2}{\left(\chi_1 + \chi_2 + 2\right)^2}, \qquad (4.15)$$

and

$$\mathcal{U}_2\left(e_1^*, e_2^*\right) = \mathcal{U}_2^* = \frac{v\left(1 + \chi_1\right)^2}{\left(\chi_1 + \chi_2 + 2\right)^2}, \qquad (4.16)$$

however, decisive for evolutionary success is only objective equilibrium profit (or equilibrium *fitness*, $\mathcal{F}_{i=1,2}$, in evolutionary terms):

$$\pi_1\left(e_1^*, e_2^*\right) = \mathcal{F}_1 = \frac{\left(\chi_1 + \chi_2 + 1\right)\left(1 + \chi_2\right) v}{\left(\chi_1 + \chi_2 + 2\right)^2}, \qquad (4.17)$$

and

$$\pi_2\left(e_1^*, e_2^*\right) = \mathcal{F}_2 = \frac{\left(\chi_1 + \chi_2 + 1\right)\left(1 + \chi_1\right) v}{\left(\chi_1 + \chi_2 + 2\right)^2}. \qquad (4.18)$$

The fitness terms determine the reproductive successes of $\chi_1$ and $\chi_2$. In the following, we use the ESS conceptions for finite and infinite populations to identify the

evolutionarily stable degree of self-confidence. As we will see, the differentiation of finite and infinite populations is mandatory in our model. Let us first consider an infinite population.

### Infinite Population

We resort to Definition 4.1 but with respect to the dimension of self-confidence. Accordingly, $\chi^{ESP}$ is an evolutionarily stable preference iff, for all $\chi \neq \chi^{ESP}$, the equilibrium condition (i) $\mathcal{F}\left(\chi^{ESP}, \chi^{ESP}\right) > \mathcal{F}\left(\chi, \chi^{ESP}\right)$, or the stability condition (ii) $\mathcal{F}\left(\chi^{ESP}, \chi^{ESP}\right) = \mathcal{F}\left(\chi, \chi^{ESP}\right)$ and $\mathcal{F}\left(\chi, \chi\right) < \mathcal{F}\left(\chi^{ESP}, \chi\right)$ hold. The following proposition shows that the players do not exhibit an evolutionarily stable disposition of over- or underconfidence in the case of $N = \infty$.

**Proposition 4.1.** For an infinite population $N = \infty$, the unique evolutionarily stable degree of self-confidence is given by

$$\chi^{ESP} = 0, \tag{4.19}$$

i.e. the players perceive their effort costs as they really are:
$c_i\left(e_i\right) = \left(1 + \chi_i^{ESP}\right) e_i = e_i$, and are profit maximizer in the sense of $\mathcal{U}_{i=1,2} = \pi_{i=1,2}$.

*Proof.* To check the conditions for ESP, consider the problem of player 1 of maximizing the equilibrium fitness function $\mathcal{F}_1 = \frac{(\chi_1 + \chi_2 + 1)(1 + \chi_2)v}{(\chi_1 + \chi_2 + 2)^2}$ (by symmetry, it is sufficient to consider only one player). The first order condition is

$$-\frac{v\left(1 + \chi_2\right)\left(\chi_1 + \chi_2\right)}{\left(\chi_1 + \chi_2 + 2\right)^3} = 0$$

After solving for $\chi_1$, which yields the best reply function

$$\chi_1\left(\chi_2\right) = -\chi_2,$$

and equating $\chi_1$ and $\chi_2$, the possible canditate for ESP is $\chi_1 = 0$.

To prove that $\chi_1 = 0$ is the only best reply against itself, consider

$$\mathcal{F}_1\left(0, 0\right) - \mathcal{F}_1\left(\chi_1, 0\right) = 1/4 \frac{v\chi_1^2}{\left(\chi_1 + 2\right)^2},$$

which is strictly positive whenever $\chi_1 \neq 0$. Thus, the equilibrium condition holds for objective payoff maximizing. $\square$

However, this does not necessarily mean that the same result appears for $N < \infty$. Let us now explore the case of a finite population.

### Finite Population

In accordance with Leininger (2009) (see also Guse and Hehenkamp, 2006), we adopt the definition of an evolutionarily stable preference (ESP) for finite populations. The transformation of Definition 4.2 to the preference frame is immediate. Accordingly, the maximization problem that the players face is of the kind as in Eq. (4.9), and it writes

$$\max_{\chi} \mathcal{F}_1(\chi, \chi^{ESS}, ..., \chi^{ESS}) - \frac{n-1}{N-1} \mathcal{F}_i(\chi, \chi^{ESS}, ..., \chi^{ESS}),$$

and for the present two-player case,

$$\max_{\chi} \mathcal{F}_1 - \frac{1}{N-1} \mathcal{F}_2, \tag{4.20}$$

and thus

$$\max_{\chi} \frac{(\chi_1 + \chi_2 + 1)(1 + \chi_2) v}{(\chi_1 + \chi_2 + 2)^2} - \frac{(\chi_1 + \chi_2 + 1)(1 + \chi_1) v}{(N-1)(\chi_1 + \chi_2 + 2)^2},$$

$$\Rightarrow \quad \max_{\chi} \frac{(\chi_1 + \chi_2 + 1) v (N - 2 + \chi_2 N - \chi_2 - \chi_1)}{(N-1)(\chi_1 + \chi_2 + 2)^2}. \tag{4.21}$$

The first order condition of this problem is given by

$$-\frac{v(2 + \chi_1 + \chi_2 + N\chi_1 + \chi_2 N\chi_1 + \chi_2 N + \chi_2^2 N)}{(N-1)(\chi_1 + \chi_2 + 2)^3} = 0. \tag{4.22}$$

Setting $\chi_1 = \chi_2 = \chi$ yields

$$-1/4 \frac{(N\chi + 1) v}{(1 + \chi)^2 (N - 1)} = 0. \tag{4.23}$$

Solving this expression for $\chi$, we reach

$$\chi = -\frac{1}{N}. \tag{4.24}$$

We have proven the following proposition.

**Proposition 4.2.** For a finite population $N < \infty$, the unique evolutionarily stable preference profile is

$$\left(\chi^{ESP}, \chi^{ESP}\right) = \left(-\frac{1}{N}, -\frac{1}{N}\right). \tag{4.25}$$

As proposition 4.2 shows, the evolutionarily stable distortion on the level of self-confidence depends on the size of the population $N$ and varies between material profit maximization for infinite populations, i.e., $\lim_{N\to\infty} \chi^{ESP} = \lim_{N\to\infty} -\frac{1}{N} = 0$, and strong overconfidence, $\chi^{ESP} = -\frac{1}{2}$, in the case of $N = 2$. The following observation is then immediate.

**Corollary 4.1.** *In any case of $2 \leq N < \infty$, the players are overconfident since $\chi^{ESP} \in [-1/2, 0)$.*

We further find the following conformity between direct evolution of effort and equilibrium effort which is induced under indirect evolution of self-confidence.

**Corollary 4.2.** *For both finite and infinite populations, and players engaged in two-player contests, direct evolution of effort and indirect evolution of self-confidence yield to the same equilibrium behavior.*

*Proof.* Substituting $\chi^{ESP} = -\frac{1}{N}$ in the equilibrium efforts (Eq. (4.14)) guides to

$$(e_1^*, e_2^*) = \left(\frac{v\left(1 + \chi_2^{ESP}\right)}{\left(\chi_1^{ESP} + \chi_2^{ESP} + 2\right)^2}, \frac{v\left(1 + \chi_1^{ESP}\right)}{\left(\chi_1^{ESP} + \chi_2^{ESP} + 2\right)^2}\right) = \left(1/4\,\frac{Nv}{N-1}, 1/4\,\frac{Nv}{N-1}\right),$$

which is the same as the ESS behavior under direct evolution, cf. section 4.2.2 or Hehenkamp et al. (2004), Theorem 5 ibidem. $\qquad \square$

## 4.4 Conclusion

Under the assumptions that evolution operates on the level of self-confidence in a finite population, we find that players overvalue their ability in two-player Tullock contests. Precisely, we detect evolutionarily stable degrees of overconfidence, meaning that the players perceive their effort costs lower than they really are. If, instead, the two-player contests appear in populations of infinite size, the players perceived efforts match exactly their true efforts, i.e. the players are objective payoff maximizer.

The fact that ESS for finite and infinite populations guides to different results is due to a simple general mechanism which holds under direct or indirect evolution. In order to succeed, evolution forces the players to maximize relative payoff. Since the players cannot affect the average payoff of all players if the population is of infinite size, the maximization problem of relative payoffs coincides with the one of absolute payoff maximization. In contrast, if the population is of finite size the players directly affect the average payoff which makes them relative payoff maximizers. This guides to the finding that the players are the more competitive, i.e. overconfidence in our setting, the smaller the population.

We further find the same equivalence result as Leininger (2009) and Boudreau and Shunda (2010) in the sense that direct evolution of effort and indirect evolution of preferences result in the same behavior. Evidently, the fact that indirect evolution of preferences and direct evolution of effort induce the same aggressive behavior under different evolving traits in Tullock contests surely deserves a more complete explanation in future research.

# Chapter 5

# Closing Remarks

Within the three main chapters of this thesis, I contributed with three essays on the evolution of individual preferences. Essay 1 and 2 are related in several aspects, however, the methodologies, purposes, and results are complementary in a wide sense. Nevertheless, with a sharp perspicacity, the results are mutually transferable to a certain extent. The third essay is scarcely related with the first and the second since it explores a complete different psychological trait in a different strategic setting.

The first essay is the study which refers most to the evolution of preferences in the literal, i.e. biological, sense. By relying on the seminal study of Bester and Güth (1998) and the extensions of Bolle (2000) and Possajennikov (2000), I tried to explain the emergence of altruistic and envious preferences in several strategic settings. The methodology to characterize the outcomes is basically biological since the replicator dynamics (where the population evolves in a polymorphic sense) and partly the adaptive dynamics (where the population evolves in a monomorphic sense) are measuring the success of the preferences. Both dynamics stemp from biology. The basic enhancement in comparison to the former approaches is that the results are dynamic, and stability is explicitly measured under the assumption of continuous preferences (i.e. the space of preferences is compact). I showed that, generally, the

qualitative findings of the former studies hold. Namely, strategic complements lead to altruism and strategic substitutes to envy.

The second essay has some assumptions which are close to those of the first one but the focal points are different. Most obviously, my purpose with the second essay is to identify evolutionary reasons why reciprocal tendencies are pervasive phenomena in economic interactions. The interpretations of the interpersonal successes are similar to those of the first essay but the perceived payoffs are more subtle. The individuals are fully identified on a two-dimensional space. In particular, they have inborn social attitudes which are not changeable by the dynamical pressures and they belong to a distribution of reciprocity parameters which is changing within the time span of a generic class of selection dynamics which works rather in a social frame of learning, imitation, or education. Hence, from a methodological view, essay 2 is a profound continuation of essay 1 in the sense that the observed time span in which the evolving trait changes is reduced from a biological to a cultural level. Another technical feature, which separates essay 2 from essay 1 somewhat, is that the methodological approach accounts for asymmetric responses in the type game. Like in essay 1, the strategic environment of the game is the fundamental determinant for measuring the type of action in the society. The results regarding reciprocity are subtle but the trend is established: it usually pays off to show a high flexibility in the reciprocal sense.

The third essay deals with individual preferences in contests. In particular, I analyzed the dimension of self-confidence in Tullock contests, where each player's winning probability is his effort's share of total efforts. I used the ESS conception for infinite populations and its analog for finite populations—while the former can rightly be seen as a refinement of Nash's concept, the latter is not necessarily one— and showed that a distorted degree of self-confidence arises only in populations of finite size; if the population under study is infinite, then only profit maximization is evolutionarily stable. Given the finite case, the players exhibit a certain degree of

overconfidence which increases with decreasing size of the population and yields to more aggressive behavior than under the undistorted Nash equilibrium level. Interestingly, the behavior which is induced by the particular degree of overconfidence under indirect evolution is the same as evolutionary stable behavior given direct evolution. Accordingly, overconfidence can be seen as another evolutionary rationale for the observation that people often exceed Nash equilibrium play in experimental contests.

In all studies the so-called indirect evolutionary approach plays a key role. Indirect evolution is currently the most established approach which provides a theoretical justification for developing subjective values in an economic society.

# Bibliography

ABRAMS, P. (2001): "Modelling the Adaptive Dynamics of Traits Involved in Inter- and Intraspecific Interactions: An Assessment of Three Methods," *Ecological Letters*, 4, 166–175.

ANDO, M. (2004): "Overconfidence in Economic Contests," Working Paper, Nihon University.

ANDREONI, J. (1990): "Impure Altruism and Donations to Public Goods: A Theory of Warm-Glow Giving," *Economic Journal*, 100, 464–477.

ANDREONI, J., P. BROWN, AND L. VESTERLUND (2002): "What Produces Fairness? Some Experimental Evidence," *Games and Economic Behavior*, 40, 1–24.

ANDREONI, J. AND J. MILLER (2002): "Giving According to GARP: An Experimental Test of the Consistency of Preferences for Altruism," *Econometrica*, 70, 737–753.

APALOO, J. (1997): "Revisiting Strategic Models of Evolution: The Concept of Neighborhood Invader Strategies," *Theoretical Population Biology*, 52, 71–77.

——— (2005): "Inaccessible Continuously Stable Strategies," *Natural Resource Modeling*, 18, 521–529.

BAHARAD, E. AND S. NITZAN (2008): "Contest Efforts in Light of Behavioural Considerations," *Economic Journal*, 118, 2047–2059.

BAYE, M., D. KOVENOCK, AND C. DE VRIES (1993): "Rigging the Lobbying Process: An Application of the All-Pay Auction," *American Economic Review*, 83, 289–294.

BECKER, G. (1976): "Altruism, Egoism, and Genetic Fitness," *Journal of Economic Literature*, 14, 817–826.

BESTER, H. AND W. GÜTH (1998): "Is Altruism Evolutionary Stable?" *Journal of Economic Behavior and Organization*, 34, 193–209.

BILLINGSLEY, P. (1968): *Convergence of Probability Measures*, Wiley, New York.

BOLLE, F. (2000): "Is Altruism Evolutionarily Stable? And Envy and Malevolence? - Remarks on Bester and Güth," *Journal of Economic Behavior and Organization*, 42, 131–133.

BOLTON, G. AND A. OCKENFELS (2000): "ERC: A Theory of Equity, Reciprocity and Competition," *American Economic Review*, 90, 166–193.

BOMZE, I. (1990): "Dynamical Aspects of Evolutionary Stability," *Monatshefte für Mathematik*, 110, 189–206.

——— (1991): "Cross Entropy Minimization in Uninvadable States of Complex Populations," *Journal of Mathematical Biology*, 30, 73–87.

BOUDREAU, J. AND N. SHUNDA (2010): "On the Evolution of Prize Perceptions in Contests," Working Paper.

BOYER, G. (1989): "Malthus Was Right After All: Poor Relief and Birth Rates in Southeastern England," *Journal of Political Economy*, 97, 93–114.

BULOW, J., J. GENEAKOPLOS, AND P. KLEMPERER (1985): "Multimarket Oligopoly: Strategic Substitutes and Complements," *Journal of Political Economy*, 93, 488–511.

CAMERER, C. (2003): *Behavioral Game Theory: Experiments on Strategic Interaction*, Princeton University Press.

CAMERER, C. AND R. THALER (1995): "More Dictator and Ultimatum Games," *Journal of Economic Perspectives*, 9, 209–219.

CHARNESS, G. AND M. RABIN (2002): "Understanding Social Preferences with Simple Tests," *Quarterly Journal of Economics*, 117, 817–869.

CHRISTIANSEN, F. (1991): "On Conditions for Evolutionary Stability for Continuously Varying Character," *American Naturalist*, 138, 37–50.

CORNES, R. AND R. HARTLEY (2003): "Loss aversion and the Tullock paradox," University of Nottingham Paper in Economics.

CRESSMAN, R. (2005): "Stability of the Replicator Equation with Continuous Strategy Space," *Mathematical Social Sciences*, 50, 127–147.

——— (2009): "Continuously Stable Strategies, Neighborhood Superiority and Two-Player Games with Continuous Strategy Space," *International Journal of Game Theory*, 38, 221–247.

CRESSMAN, R. AND J. HOFBAUER (2005): "Measure Dynamics on a One-Dimensional Continuous Trait Space: Theoretical Foundations for Adaptive Dynamics," *Journal of Theoretical Biology*, 67, 47–59.

CRESSMAN, R., J. HOFBAUER, AND F. RIEDEL (2006): "Stability of the Replicator Equation for a Single Species with a Multi-Dimensional Trait Space," *Journal of Theoretical Biology*, 239, 273–288.

DEBREU, G. (1959): *Theory of Value: An Axiomatic Analysis of Economic Equilibrium*, Wiley, New York.

DEKEL, E., J. ELY, AND O. YILANKAYA (2007): "Evolution of Preferences," *Review of Economic Studies*, 74, 685–704.

DOHMEN, T. AND A. FALK (2006): "Performance Pay and Multi-Dimensional Sorting: Productivity, Preferences and Gender," IZA Discussion Paper, No. 2001.

DUFWENBERG, M. AND G. KIRCHSTEIGER (2004): "A Theory of Sequential Reciprocity," *Games and Economic Behavior*, 47, 268–298.

ECKEL, C. AND P. GROSSMAN (1998): "Are woman less selfish than men? Evidence from dictator experiments," *Economic Journal*, 108, 726–735.

ELY, J. AND O. YILANKAYA (2001): "Nash Equilibrium and the Evolution of Preferences," *Journal of Economic Theory*, 97, 255–272.

ESHEL, I. (1983): "Evolutionary and Continuous Stability," *Journal of Theoretical Biology*, 103, 99–111.

ESHEL, I. AND U. MOTRO (1981): "Kin Selection and Strong Evolutionary Stability of Mutual Help," *Theoretical Population Biology*, 19, 420–433.

ESHEL, I. AND E. SANSONE (2003): "Evolutionary and Dynamic Stability in Continuous Population Games," *Journal of Mathematical Biology*, 46, 445–459.

FALK, A. AND U. FISCHBACHER (2006): "A Theory of Reciprocity," *Games and Economic Behavior*, 54, 293–315.

FEHR, E. AND U. FISCHBACHER (2002): "Why Social Preferences Matter–the Impact of Non-selfish Motives on Competition, Cooperation, and Incentives," *The Economic Journal*, 112, C1–C33.

——— (2003): "The Nature of Human Altruism," *Nature*, 425, 785–791.

FEHR, E. AND K. SCHMIDT (1999): "A Theory of Fairness, Competition, and Cooperation," *Quarterly Journal of Economics*, 114, 817–868.

FERSHTMAN, C. AND Y. WEISS (1998): "Social Rewards, Externalities and Stable Preferences," *Journal of Public Economics*, 70, 53–70.

FRANK, R. (1987): "If Homo Economicus Could Choose His Own Utility Function, Would He Want One with a Conscience?" *American Economic Review*, 77, 593–604.

——— (1988): *Passions within Reason: The strategic Role of the Emotions*, W.W. Norton, New York.

GEANAKOPLOS, J., D. PEARCE, AND E. STACCHETTI (1989): "Psychological Games and Sequential Rationality," *Games and Economic Behavior*, 1, 60–79.

GIGERENZER, G. AND R. SELTEN (2001): *Bounded rationality: The adaptive toolbox*, Cambridge, MA: MIT Press.

GINTIS, H. (2000): "Strong Reciprocity and Human Sociality," *Journal of Theoretical Biology*, 206, 169–179.

GÜTH, S., W. GÜTH, AND H. KLIEMT (2002): "The Dynamics of Trustworthiness Among the Few," *Japanese Economic Review*, 53, 369–388.

GÜTH, W. AND A. OCKENFELS (2005): "The Coevolution of Morality and Legal Institutions: An Indirect Evolutionary Approach," *Journal of Institutional Economics*, 1, 155–174.

GÜTH, W. AND B. PELEG (2001): "When Will Payoff Maximization Survive? An Indirect Evolutionary Analysis," *Journal of Evolutionary Economics*, 11, 479–499.

GÜTH, W., R. SCHMITTBERGER, AND B. SCHWARZE (1982): "An Experimental Analysis of Ultimatum Bargaining," *Journal of Economic Behavior and Organisation*, 3, 367–388.

GÜTH, W. AND M. YAARI (1992): "Explaining Reciprocal Behavior in Simple Strategic Games: An Evolutionary Approach," in *Explaining Forces and Changes: Approaches to Evolutionary Economics*, ed. by U. Witt., University of Michigan Press.

GUSE, T. AND B. HEHENKAMP (2006): "The Strategic Advantage of Interdependent Preferences in Rent-Seeking Contests," *Public Choice*, 129, 323–352.

GUTTMAN, J. (2000): "On the Evolutionary Stability of Preferences for Reciprocity," *European Journal of Political Economy*, 16, 31–50.

HAMILTON, W. (1971): "Selfish and spiteful behavior in an evolutionary model," *Nature*, 228, 1218–1220.

HARRISON, R. AND M. VILLENA (2008): "On the Evolution of Reciprocal Behavior: A Game Theoretic Approach," Working Paper.

HEHENKAMP, B., W. LEININGER, AND A. POSSAJENNIKOV (2004): "Evolutionary Equilibrium in Tullock Contests: Spite and Overdissipation," *European Journal of Political Economy*, 20, 1045–1057.

HEIDHUES, P. AND F. RIEDEL (2007): "Do Social Preferences Matter in Competitive Markets?" Working Paper, Bielefeld University.

HEIFETZ, A. AND E. SEGEV (2004): "The Evolutionary Role of Toughness in Bargaining," *Games and Economic Behaviour*, 49, 117–134.

HEIFETZ, A., C. SHANNON, AND Y. SPIEGEL (2007a): "The Dynamic Evolution of Preferences," *Economic Theory*, 32, 251–286.

——— (2007b): "What to Maximize if You Must," *Journal of Economic Theory*, 133, 31–57.

HOFBAUER, J., J. OECHSSLER, AND F. RIEDEL (2009): "Brown-von Neumann-Nash Dynamics: The Continuous Strategy Case," *Games and Economic Behavior*, 65, 406–429.

HOFBAUER, J., P. SCHUSTER, AND K. SIGMUND (1979): "A Note on Evolutionary Stable Strategies and Game Dynamics," *Journal of Theoretical Biology*, 81, 609–612.

HOFFMAN, E., K. MCCABE, AND V. SMITH (1996): "Social distance and other-regarding behavior in dictator games," *American Economic Review*, 86, 653–660.

HOFFMANN, E., K. MCCABE, AND V. SMITH (1998): "Behavioral Foundations of Reciprocity: Experimental Economics and Evolutionary Psychology," *Economic Inquiry*, 36, 335–352.

HÖRISCH, H. AND O. KIRCHKAMP (2010): "Less Fighting than Expected: Experiments with Wars of Attrition and All-Pay Auctions," *Public Choice*, 144, 347–367.

HUCK, S. AND J. OECHSSLER (1999): "The Indirect Evolutionary Approach to Explaining Fair Allocations," *Games and Economic Behavior*, 28, 13–24.

KISDI, É. AND G. MESZÉNA (1995): "Life Histories with Lottery Competition in a Stochastic Environment: ESSs which do not Prevail," *Theoretical Population Biology*, 47, 191–211.

KONRAD, K. (2004): "Altruism and Envy in Contests: An Evolutionarily Stable Symbiosis," *Social Choice and Welfare*, 22, 479–490.

KONRAD, K. AND H. SCHLESINGER (1997): "Risk aversion in rent seeking and rent augmenting games," *Economic Journal*, 107, 1671–1683.

KYLE, A. AND A. WANG (1997): "Speculation Duopoly with Agreement to Disagree: Can Overconfidence Survive the Market Test?" *The Journal of Finance*, 52, 2073–2090.

LEININGER, W. (2009): "Evolutionarily stable preferences in contests," *Public Choice*, 140, 341–356.

LEVINE, J. (1998): "Modeling Altruism and Spitefulness in Experiments," *Review of Economic Dynamics*, 1, 593–622.

LUDWIG, S., P. WICHARDT, AND H. WICKHORST (2010): "Overconfidence Can Improve an Agent's Relative and Absolute Performance in Contests," Working Paper.

MARROW, P., U. DIECKMANN, AND R. LAW (1996): "Evolutionary Dynamics of Predator-Prey Systems: An Ecological Perspective," *Journal of Mathematical Biology*, 34, 556–578.

MAYNARD-SMITH, J. (1982): *Evolution and the Theory of Games*, Cambridge University Press.

MAYNARD-SMITH, J. AND G. PRICE (1973): "The Logic of Animal Conflicts," *Nature*, 246, 15–18.

MORGAN, J., H. ORZEN, AND M. SEFTON (2008): "Endogenous Entry in Contests," Working Paper, University of Nottingham.

MOULIN, H. (1984): "Dominance Solvability and Cournot Stability," *Mathematical Social Sciences*, 7, 83–102.

MUI, V. (1995): "The Economics of Envy," *Journal of Economic Behavior and Organisation*, 26, 311–336.

NALEBUFF, B. AND J. STIGLITZ (1983): "Prizes and Incentives: Towards a General Theory of Compensation and Competition," *Bell Journal of Economics*, 14, 21–43.

NASH, J. (1950): *Non-cooperative Games*, Ph.D. dissertation, Princeton University.

NOWAK, M. AND K. SIGMUND (2005): "Evolution of Indirect Reciprocity," *Nature*, 437, 1291–1298.

OECHSSLER, J. AND F. RIEDEL (2001): "Evolutionary Dynamics on Infinite Strategy Spaces," *Economic Theory*, 17, 141–162.

——— (2002): "On the Dynamic Foundation of Evolutionary Stability in Continuous Models," *Journal of Economic Theory*, 107, 223–252.

OK, E. AND F. VEGA-REDONDO (2001): "On the Evolution of Individualistic Preferences: An Incomplete Information Scenario," *Journal of Economic Theory*, 97, 231–254.

POSSAJENNIKOV, A. (1999): "On Evolutionary Stability of Spiteful Preferences," Working Paper, CentER Tilburg University.

——— (2000): "On the Evolutionary Stability of Altruistic and Spiteful Preferences," *Journal of Economic Behavior and Organisation*, 42, 125–129.

RABIN, M. (1993): "Incorporating Fairness into Game Theory and Economics," *American Economic Review*, 83, 1281–1302.

ROBSON, A. AND S. SKAPERDAS (2008): "Costly Enforcement of Property Rights and the Coase Theorem," *Economic Theory*, 36, 109–128.

ROTH, A. (1995): "Bargaining Experiments," in *Handbook of Experimental Economics*, ed. by J. Kagel and A. Roth, Princeton University Press, Princeton.

SAMUELSON, L. AND J. ZHANG (1992): "Evolutionary Stability in Asymmetric Games," *Journal of Economic Theory*, 57, 363–391.

SCHAFFER, M. (1988): "Evolutionarily Stable Strategies for a finite Population and a variable Contest Size," *Journal of Theoretical Biology*, 132, 469–478.

SCHELLING, T. (1960): *The Strategy of Conflict*, Harvard University Press, Cambridge.

SCHMIDT, F. (2009): "Evolutionarily Stability of Altruism and Envy in Tullock Contests," *Economic Governance*, 10, 247–259.

SELTEN, R. (1980): "A Note on Evolutionary Stable Strategies in Asymmetric Animal Conflicts," *Journal of Theoretical Biology*, 84, 93–101.

——— (1991): "Evolution, Learning, and Economic Behavior," *Games and Economic Behavior*, 3, 3–24.

SETHI, R. AND E. SOMANTHAN (2001): "Preference Evolution and Reciprocity," *Journal of Economic Theory*, 97, 273–297.

——— (2003): "Understanding Reciprocity," *Journal of Economic Behavior and Organization*, 50, 1–27.

SHIRYAEV, A. (1995): *Probability*, Springer, Berlin.

SOBEL, J. (2005): "Interdependent Preferences and Reciprocity," *Journal of Economic Literature*, 43, 396–440.

SZYMANSKI, S. (2003): "The Economic Design of Sporting Contests," *Journal of Economic Literature*, 41, 1137–1187.

TAYLOR, P. (1989): "Evolutionary Stable Strategies and Game Dynamics," *Theoretical Population Biology*, 36, 125–143.

TAYLOR, P. AND L. JONKER (1978): "Evolutionary Stable Strategies and Game Dynamics," *Mathematical Biosciences*, 40, 145–156.

TRIVERS, R. (1971): "The Evolution of Reciprocal Altruism," *Quarterly Review of Biology*, 46, 35–58.

TULLOCK, G. (1980): "Efficient rent-seeking," in *Toward a Theory of the Rent-Seeking Society*, ed. by Buchanan/Tollison/Tullock, College Station: Texas A and M University Press.

VON NEUMANN, J. AND O. MORGENSTERN (1944): *Theory of Games and Economic Behvavior*, Princeton University Press.

WEIBULL, J. (1995): *Evolutionary Game Theory*, The MIT Press, Cambridge, MA.

# Kurzer Lebenslauf

Niko Noeske

Geboren am 27.02.1978 in Bonn

1998  Abitur am Pädagogium Otto-Kühne Schule, Bonn

1999  Zivildienst, Untere Landschaftsbehörde, Bonn

Studium der Volkswirtschaftslehre an der Rheinischen-Friedrich-Wilhelms-Universität Bonn

2006  Diplom-Volkswirt

Promotionsstudium an der Universität Bielefeld

2011  Promotion zum Dr. rer. pol. an der Wirtschaftswissenschaftlichen Fakultät, Universität Bielefeld