

ECONOMIC STUDIES
DEPARTMENT OF ECONOMICS
SCHOOL OF BUSINESS, ECONOMICS AND LAW
UNIVERSITY OF GOTHENBURG
173

Essays on Epistemology and Evolutionary Game Theory

Elias Tsakas

ISBN 91-85169-32-3
ISBN 978-91-85169-32-0
ISSN 1651-4289 print
ISSN 1651-4297 online



UNIVERSITY OF GOTHENBURG

ELIAS TSAKAS

Essays on Epistemology and Evolutionary
Game Theory

May, 2008

Department of Economics, Göteborg University

ISSN:

to the memory of Jolene

Preface

Despite what people believe, the process for receiving a PhD is quite simple: you come across many ideas, you get confused, you try to understand, you get more confused, you try harder and hopefully in the end you do not manage to completely clarify things in your mind. If you still feel confused after you become a doctor, then you have made a step towards becoming an academic, and you are ready to conduct research for the rest of your life. During this struggle I met many really interesting people who “contributed to my confusion”, and for that I am really grateful to them. Although the space is limited I will try to pay a tribute to all of them.

I would like to start by thanking my thesis advisor/co-author/friend Mark Voorneveld, who has done a superb job in every aspect. Not only has he taught me how to write papers, but he has also been very supportive and incredibly patient. He opened his home to me and I really enjoyed having lunch or dinner with him and his family, Fia and Femke, whenever I was in Stockholm. I truly believe that without his help I would have never been able to complete this thesis. I really hope that we will be in touch for the years to come and we will continue working together.

While Mark was my official advisor, I was really lucky to have a number of prominent game theorists around me, who in many instances acted like informal advisors and I owe them a great deal of grace for that.

Martin Dufwenberg has been visiting Göteborg University on a regular basis during the last three years. During these visits he spent so much time reading my papers and talking about them. Martin’s help has been really invaluable in every aspect and I feel truly lucky to have met him.

When I decided to have an external thesis advisor, it seemed as if I was taking a huge risk. Olof Johansson-Stenman, the vice-head of our department and director of the PhD program was the one who supported my decision from the beginning to the end. He acted as my internal advisor the entire time, he was the one I would go and talk whenever a problem of academic nature popped up and for that I honestly feel the need to thank him.

The third chapter of this thesis is based on a paper that I wrote jointly with Olivier Gossner. I met Olivier in the annual game theory conference organized by the State University of New York at Stony Brook and we decided to work together towards article. I feel really lucky to have worked with him since he is widely considered as one of the greatest game theorists of his generation. During our collaboration I learned a lot on how to do research for that I would like to thank Olivier.

David Ahn was appointed my thesis advisor during the time I spent in the University of California, Berkeley. Usually this kind of appointments do not mean much: you meet once in the beginning of the academic year and never again. However, that was not the case this time. David really made me feel part of the department. His door was always open to me and I owe him a lot for that.

During my stay at Berkeley I met many amazing people. The one person I would like to specially mention is Amanda Friedenberg, who happened to be visiting Berkeley at the same time as me. Amanda is one of the few people who are interested in game theoretic models of epistemology. Since she is considered as one of the brightest new theorists I feel lucky to have worked together with her. I would really like to thank Amanda for the fact that she spent so much time talking with me about research issues, and also for having invited me to visit her home department at the Olin Business School, at the Washington University in St. Louis.

I would like to thank everybody who took time out of their work to read chapters of this thesis and give very useful comments which ultimately helped me improve the quality of the thesis. I would like to start by thanking Geir Asheim, who accepted to act as the main opponent in my dissertation defense and the members of my PhD committee, Jörgen Weibull, Martin Dufwenberg and Jeff Steif. I would also like to specially thank Andrés Perea, who acted as the discussant in the final seminar of my PhD the-

sis. The points brought up by David Ahn, Alpaslan Akay, Giacomo Bonanno, Martin Dufwenberg, Amanda Friedenberg, John Geanakoplos, Aviad Heifetz, Barton Lipman, Lucie Menager, Andrew Postlewaite, Fernando Vega-Redondo, Dov Samet, Bill Sandholm, Jörgen Weibull, the seminar attendants at UC Berkeley, Göteborg University, Maastricht University, Cardiff Business School and the participants in the 16th and 17th International Conference on Game Theory at the State University of New York in Stony Brook, the 6th Conference for Research in Economic Theory and Econometrics (CRETE) in Naxos, Greece, the Conference in Honor of Vernon Smith at the University of Arizona in Tucson and the 3rd Royal Economic Society PhD meeting at UC London have been of great help. Their contribution has been really important and is greatly appreciated.

An integral part of my graduates studies has been the courses I attended during the first two years of the program. I believe that without this coursework I would not have been able to complete my dissertation. For that I am really grateful to my teachers: Jörgen Weibull, Mark Voorneveld, Olof Johansson-Stenman, Fredrik Carlsson, Peter Martinson, Douglas Hibbs, Lennart Flood, Lennart Hjalmarsson, Ola Olsson, Renato Aguilar, Katarina Nordblom, Ali Tasiran, Roger Wahlberg, Henry Olsson, Bo Sandelin, Karl-Markus Modén, Hong Wu and Steffen Jørgensen.

I would like to gratefully acknowledge the help of the administration staff for their contribution to the completion of this thesis. I would like to particularly thank Eva-Lena Neth-Johansson, who was the first person I would go to whenever I faced a practical problem. If Eva-Lena does not work, nothing works in our department. I would also like to thank Eva Jonasson, Elizabeth Földi, Katarina Renström Jeanette Saldjoughi, and Göran Persson for their invaluable help at various administrative issues throughout my doctoral studies in Göteborg. Finally, I would like to thank Monica Allen from UC Berkeley, Ritva Kiviharju from the Stockholm School of Economics and Nanci Kirk from the Olin Business School at Washington University in St. Louis for making my life easier during my stay there.

Large part of dissertation was written during my academic visits. I am particularly grateful to the Department of Economics at the University of California, Berkeley, where I spent a considerable part of my PhD studies. I would like to particularly thank the microeconomic theory group at Berkeley for warmly welcoming me and making me feel

a part of it from the first moment. I would also like to thank the Stockholm School of Economics for its hospitality every time I visited Mark, and the Olin Business School at Washington University in St. Louis, that I visited towards the end of my studies, invited by Amanda Friedenbergl. In addition, I would like to thank the Economics Departments of the Stockholm School of Economics, Aarhus University and Växjö University for letting me take courses offered in their PhD program. Finally, I would like to mention that part of this thesis was written in the Doma cafe in New York City, the Free Speech cafe in the University of California at Berkeley and the Nöller espressobar in Göteborg, which provided a really nice and stimulating environment for thinking.

I would like to gratefully acknowledge the financial support from the Adlerbertska Forskningsstiftelsen, the Adlerbertska Stipendiestiftelsen, the Adlerbertska Hospitiestiftelsen and the Stiftelsen Stipendiefonden Viktor Rydbergs Minne. I would also like to thank the South Swedish Graduate Program in Economics, the Nordic Network in Economics, the Stiftelsen Henrik Ahrenbergs Studiefond and the Stiftelsen Paul och Marie Berghaus Donationsfond for financing my traveling costs to conferences and part of my academic visits to other departments.

I would like to thank my classmates and friends Andreea Mitrut, Bruno Turnheim, Florin Maican, Marcela Ibañez, Martine Visser, Jorge Garcia, Conny Wollbrant, Jonas Alem, Clara Villegas, Anna Widerberg, Sven Tengstam, Jiegen Wei, Qin Ping, Precious Zhikali, Daniel Zerfu, Miguel Quiroga, Innocent Kabenga, Gustav Hansson, Elina Lampi, Annika Lindskog, Ann-Sofie Isaksson, Pelle Ahlerup, Niklas Jakobsson, Miyase Köksal, Anton Nivorozhkin, Violeta Piculescu and Kerem Tezic. Knowing that they were going through the same process made me feel that I was not alone in this. I would also like to thank Johan Lönnroth and Wlodek Bursztyn, with who I shared many great conversations over politics and other issues of general interest during my breaks.

My Greek friends in Göteborg, Christos Divanis, Giannis Basinos, Giannis Milionis, Kostas Eleftherakis, Ellen Lekka, Alex Balatsos, Andreas Vandoros, Nikos Papachatzakis, Giannis Lennart Balatsos and Aris Seferiadis were always next to me during these years. The moments of laugh and joy that we shared where really vital for my me in order to continue and complete my doctoral studies. Very important for the completion of my thesis has been my friendship with Stefan Heldmann, Casper Holmgaard Jensen and the

rest of my roommates in the International House in Berkeley. Invaluable was the support of all my friends who were away but with whom I regularly spoke on the phone during my doctoral studies.

I would like to specially thank my best friend in Göteborg Alpaslan Akay, with whom I spent many hours talking not only about econometrics, statistics and game theory, but also about Wittgenstein, Karl Popper and Pink Floyd, among other things. My first publication is joint work with Aslan and I truly hope there will be many more in the future.

Last but certainly not least, I would like to thank my family who has always been there for me. My father Thomas and my mother Roula have always believed in me, which has made me confident that I can achieve anything I wish. Finally, I would like to thank my brother Nikolas, who is the most important person in my life, for being always on my side supporting every decision I take. I honestly do not think I would have been able to achieve anything without them.

Göteborg, May 2008

Elias Tsakas

Abstract

This thesis has two parts, one consisting of three independent papers in epistemology (Chapters 1-3) and another one consisting of a single paper in evolutionary game theory (Chapter 4):

- (1) “Knowing who speaks when: A note on communication, common knowledge and consensus” (together with Mark Voorneveld)

We study a model of pairwise communication in a finite population of Bayesian agents. We show that, if the individuals update only according to the signal they actually hear, and they do not take into account all the hypothetical signals they could have received, a consensus is not necessarily reached. We show that a consensus is achieved for a class of protocols satisfying “information exchange”: if agent A talks to agent B infinitely often, agent B also gets infinitely many opportunities to talk back. Finally, we show that a commonly known consensus is reached in arbitrary protocols, if the communication structure is commonly known.

- (2) “Aggregate information, common knowledge and agreeing not to bet”

I consider gambles that take place even if some – but not all – people agree to participate. I show that the bet cannot take place if it is commonly known how many individuals are willing to participate.

- (3) “Testing rationality on primitive knowledge” (together with Olivier Gossner)

The main difficulty in testing negative introspection is the infinite cardinality of the set of propositions. We show that, under positive conditions, negative introspection holds if and only if it holds for primitive propositions, and is therefore

easily testable. When knowledge arises from a semantic model, we show that, further, negative introspection on primitive propositions is equivalent to partitional information structures. In this case, partitional information structures are easily testable.

(4) “The target projection dynamic” (together with Mark Voorneveld)

We study a model of learning in normal form games. The dynamic is given a microeconomic foundation in terms of myopic optimization under control costs due to a certain status-quo bias. We establish a number of desirable properties of the dynamic: existence, uniqueness, and continuity of solution trajectories, Nash stationarity, positive correlation with payoffs, and innovation. Sufficient conditions are provided under which strictly dominated strategies are wiped out. Finally, some stability results are provided for special classes of games.

KEYWORDS: Common knowledge, communication, consensus, betting, primitive propositions, negative introspection, information partition, projection, learning.

JEL CODES: C72, D80, D81, D82, D83, D84, D89.

Contents

Preface V

Abstract XI

Part I Epistemology

1 Communication, common knowledge and consensus 3

1.1 Introduction 3

1.2 Notation and preliminaries 5

1.2.1 Information and knowledge 5

1.2.2 Signals 5

1.2.3 Communication protocol 6

1.3 Actual and hypothetical signals 7

1.4 Common knowledge of the protocol and consensus 9

Appendix 12

2 Agreeing not to bet 15

2.1 Introduction 15

2.2 Knowing how many players participate 17

2.2.1 Information and knowledge 17

2.2.2 Gambles with limited participation 17

2.2.3 Main result 18

Appendix 19

3	Testing rationality on primitive knowledge	23
3.1	Introduction	23
3.2	Knowledge	25
3.3	Primitive propositions and negative introspection	27
3.4	Primitive propositions in semantic models of knowledge	30
	Appendix	32

Part II Evolutionary Game Theory

4	The target projection dynamic	35
4.1	Introduction	35
4.2	Notation and preliminaries	37
4.2.1	Learning in normal form games	37
4.2.2	Projections	38
4.3	The target projection dynamic	39
4.3.1	General properties	42
4.3.2	Strict domination: mind the gap	44
4.3.3	The projection dynamic and the target projection dynamic	46
4.4	Special classes of games	48
4.4.1	Stable games	48
4.4.2	Zero-sum games	49
4.4.3	Games with strict Nash equilibria	49
4.4.4	Games with evolutionarily stable strategies	50
	Appendix	52
	References	55

Epistemology

Knowing who speaks when: A note on communication, common knowledge and consensus

1.1 Introduction

Aumann (1976) showed in his seminal paper that if two people have the same prior, and their posteriors for an event are common knowledge, then these posteriors are identical. Geanakoplos and Polemarchakis (1982) put this result in a dynamic framework by showing that if they communicate their posteriors back and forth, they will eventually agree on a common probability assessment. Cave (1983) and Bacharach (1985) independently generalized these results to finite populations and arbitrary signal functions, in place of posterior probabilities. Their setting has been the stepping stone for further development of models of communication in populations with Bayesian agents. The main aim of this literature is to study the conditions for reaching a consensus in groups of people through different communication mechanisms.

All the previous models assume that communication takes place through public announcement of the signals, which is quite restrictive. Parikh and Krasucki (1990) relaxed this assumption by introducing a model of pairwise private communication. They showed that under some connectedness assumption on the structure of the communication protocol (fairness), i.e., if everybody talks to everybody – directly or indirectly – consensus will be reached. A number of subsequent papers studied consensus in environments with pairwise communication, under different assumptions about the signal functions, the protocol structure and information structure (see Krasucki, 1996; Heifetz, 1996; Koessler, 2001; Houy and Menager, 2007).

A common assumption in all models with pairwise communication is that when an individual receives a signal she does not condition only on what she hears, but she takes into account all different hypothetical scenarios that could have occurred, had the sender of the signal acted differently. Thus, the recipient implicitly processes much more information than the one embodied in the actual signal.

In the present paper we relax this assumption. Instead we suppose that whenever an individual receives a signal, she conditions on what she hears, and not on all contingent scenarios. That is, individuals take into account only the actual signals, and not the hypothetical ones. As we show, taking into account only the actual signals does not suffice for consensus when the population communicates through an arbitrary fair protocol. A partial result can be established instead: if individuals update their information given only the actual signals and the communication protocol satisfies information exchange, i.e., if i cannot talk to j without hearing from j , then a consensus is reached.

In the second part of the paper we provide sufficient conditions for consensus through an arbitrary fair protocol. Assume that the structure of the protocol is commonly known. In this a case, individuals who do not participate in the conversation (third parties) learn something from their knowledge about the structure of the protocol. The information that third parties receive is summarized in the set of states that are consistent with the idea of the sender having talked to the receiver. That is, third parties rule out the signals that the receiver could not have heard, and condition on the rest. Clearly, since the true signal cannot be ruled out, third parties condition on a larger set, implying that what they learn – from their knowledge about the protocol – is less informative than the actual signal that the receiver hears. Then, we show that consensus will be achieved if communication takes place according to any fair protocol. This follows from the fact that common knowledge of the protocol, induces common knowledge of the signals.

The previous result is quite surprising since not everybody hears the signals, and not everybody can see that everybody hears the signal, and so on, as it would happen with public announcement of the signals. To see this consider the following example, due to Heifetz (1996). Alice, Bob and Carol sit in a circle. Alice observes the outcome of tossing a fair coin, and whispers it in Bob's ear. When Carol sees Alice talking to Bob, she does not know what she has told him, but she knows that they have spoken. In such a structure

common knowledge is not a natural consequence of the announcement, and therefore the result is not obvious.

1.2 Notation and preliminaries

1.2.1 Information and knowledge

Consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and a finite population $N = \{1, \dots, n\}$. The measure \mathbb{P} determines the common prior beliefs of the individuals in the population about every event $E \in \mathcal{F}$. Every individual is endowed with a finite **information partition** $\mathcal{I}_i \subseteq \mathcal{F}$. Let $I_i(\omega) \in \mathcal{I}_i$ contain the states that i cannot distinguish from ω . In other words $I_i(\omega)$ denotes i 's private information at ω . We say that i satisfies non-delusion if she does not rule out the true state of the world, i.e., if $\omega \in I_i(\omega)$ for all $\omega \in \Omega$. Throughout the paper we assume non-delusion.

Let $\mathcal{J} = \vee_{i=1}^n \mathcal{I}_i$, and $\mathcal{M} = \wedge_{i=1}^n \mathcal{I}_i$ denote the join (coarsest common refinement), and the meet (finest common coarsening) of the information partitions respectively. Similarly to Geanakoplos and Polemarchakis (1982), we assume¹ that $\mathbb{P}[J] > 0$ for every $J \in \mathcal{J}$. We define knowledge as usual, i.e., we say that i **knows** some $E \in \mathcal{F}$ at ω if $I_i(\omega) \subseteq E$. The event E is **commonly known** at ω if $M(\omega) \subseteq E$, where $M(\omega)$ denotes the member of the meet \mathcal{M} that contains ω .

1.2.2 Signals

Let A be a non-empty set of actions. A **signal (action) function** $f_i : \Omega \rightarrow A$ determines what signal agent i transmits at every $\omega \in \Omega$. We assume that an individual i transmits the same signal at every state of her information set, i.e., $f_i(\omega') = f_i(\omega)$ for every $\omega' \in I_i(\omega)$ and every $\omega \in \Omega$. Individual i 's signal is commonly known at some state ω , if $M(\omega) \subseteq R_i$, where

$$R_i = \{\omega' \in \Omega : f_i(\omega') = f_i(\omega)\}. \quad (1.1)$$

Consensus has been reached at some state ω if all individuals transmit the same signal while being at ω , i.e., if $f_i(\omega) = f_j(\omega)$ for all $i, j \in N$.

¹ Brandenburger and Dekel (1987) relax this assumption.

Let the *working partition* (Heifetz, 1996; Krasucki, 1996) be the coarsening \mathcal{P}_i of \mathcal{I}_i defined as follows: two arbitrary states $\omega, \omega' \in \Omega$ lie in the same element of \mathcal{P}_i if and only if $f_i(\omega) = f_i(\omega')$. In other words, \mathcal{P}_i is the collection of individual i 's equivalent classes of messages. By definition, $R_i \in \mathcal{P}_i$. Let $P_i(\omega)$ denote the element of \mathcal{P}_i that contains the state ω .

Let $\sigma(\mathcal{J})$ be the σ -algebra generated by \mathcal{J} . Agents in the population are *like-minded* if there is a function $f : \sigma(\mathcal{J}) \rightarrow A$, called the virtual signal function, such that $f_i(\omega) = f(I_i(\omega))$ for every $i \in N$ and $\omega \in \Omega$. A function f satisfies *union consistency*² (Cave, 1983) if for all disjoint $J_1, J_2 \in \sigma(\mathcal{J})$ with $f(J_1) = f(J_2)$, it holds that $f(J_1 \cup J_2) = f(J_1)$. Henceforth, we assume that f is real-valued. If the signals are posterior beliefs about some event E , the function f_i can be rewritten as $f_i(\omega) = \mathbb{P}[E|I_i(\omega)]$. Parikh and Krasucki (1990) considered the following stronger version of union consistency: a signal function $f : \sigma(\mathcal{J}) \rightarrow \mathbb{R}$ satisfies *convexity*, if for all non-empty, disjoint $J_1, J_2 \in \sigma(\mathcal{J})$ there is some $\alpha \in (0, 1)$ such that $f(J_1 \cup J_2) = \alpha f(J_1) + (1 - \alpha)f(J_2)$.

1.2.3 Communication protocol

Let N be a population of like-minded individuals. Every individual $i \in N$ is endowed with an \mathcal{F} -measurable information partition \mathcal{I}_i^1 at time $t = 1$, i.e., before any communication takes place. At every $t \in \mathbb{N}$ a sender s_t privately announces her signal to a recipient r_t , who updates her information to $\mathcal{I}_{r_t}^{t+1}$ according to some refining mechanism. Individuals $j \neq r_t$ do not receive any signals and consequently do not revise their information. Communication then proceeds to the next stage $t + 1$. The communication pattern, i.e., who talks to whom at each period is determined by the sequence $\{(s_t, r_t)\}_{t \in \mathbb{N}}$ in $N \times N$, referred to as the *protocol*.

The protocol induces a graph on N with a directed edge from i to j , if i talks to j infinitely often, i.e., if there are infinitely many $t \in \mathbb{N}$ with $(s_t, r_t) = (i, j)$. Parikh and Krasucki (1990) called a protocol *fair* if the graph of directed edges is strongly connected, i.e., if there is a path of directed edges which starts from some individual, passes from all the vertexes (individuals), returning to its origin. In other words, a protocol is fair if for all distinct $i, j \in N$, there is a path from i to j and back.

² Bacharach (1985) used the term *sure-thing principle* for the same property.

A protocol satisfies *information exchange* if for all distinct $i, j \in N$ with a directed edge from i to j , there is a directed edge from j to i (Krasucki, 1996).

1.3 Actual and hypothetical signals

Parikh and Krasucki (1990) were the first ones to study the conditions under which a population reaches a consensus through pairwise communication. They showed that if the protocol is fair and the signals convex, a consensus will be eventually achieved³. Krasucki (1996) consequently proved that this result can be generalized to union consistent signals, for a special class of protocols, i.e., those which satisfy information exchange.

However, if we take a closer look at the mechanism that the recipients use when they update their information, we will see that they actually take into account not only the actual signal transmitted by the sender, but all hypothetical scenarios that could have occurred, had s_t sent some other signal. In other words, the receiver implicitly considers what would have happened, had the sender said something else and then refines her partition accordingly. Then she repeats this process for every possible signal.

Formally, the information partitions are refined as follows at time t :

$$\mathcal{I}_j^{t+1} = \begin{cases} \mathcal{I}_j^t & \text{if } j \neq r_t, \\ \mathcal{I}_j^t \vee \mathcal{P}_{s_t}^t & \text{if } j = r_t. \end{cases} \quad (1.2)$$

The receiver scans in his mind the entire state space, and at every state $\omega' \in \Omega$ she conditions on the message $f_{s_t}^t(\omega')$ that would have been announced by the sender, had the true state been ω' . This mechanism was introduced⁴ by Parikh and Krasucki (1990), and was adopted by all subsequent papers in the literature (Parikh, 1996; Heifetz, 1996; Koessler, 2001; Houy and Menager, 2007).

Though one can find examples involving communication where the agents can contemplate where the signals come from, and therefore run all the hypothetical scenarios in their mind, this is not always the case. Quite often the recipient gets to improve her information

³ It can be shown actually that, contrary to what they argue, this consensus will be commonly known.

⁴ Parikh and Krasucki (1990) actually use a slightly different recursive definition of the updating process, which however can be easily shown that it is equivalent to (1.2). The recursive version of (1.2) was also used by Koessler (2001), and Houy and Menager (2007).

by taking into account only the actual signal, rather than all possible contingencies. In this case the refining mechanism at time t becomes

$$\mathcal{I}_j^{t+1} = \begin{cases} \mathcal{I}_j^t & \text{if } j \neq r_t, \\ \mathcal{I}_j^t \vee \mathcal{R}_{s_t}^t & \text{if } j = r_t, \end{cases} \quad (1.3)$$

where $\mathcal{R}_{s_t}^t = \{R_{s_t}^t, (R_{s_t}^t)^c\}$. Clearly when the individuals refine their partitions according to (1.3) they process less information than when they use (1.2), since $\mathcal{R}_{s_t}^t$ is coarser than $\mathcal{P}_{s_t}^t$. Therefore, any result assuming (1.3) instead of (1.2) is stronger.

Finally, refining according to (1.3) implies that the individuals actually communicate, i.e., s_t says something to r_t at time t and r_t updates her information given what she has heard. Using (1.2) on the other hand, implies hypothetical communication. That is, r_t learns that s_t is about to talk to her at time t , takes a look at $\mathcal{I}_{s_t}^t$ and updates her information, before even having heard the signal $f_{s_t}^t(\omega)$.

The first natural question that arises at this point is whether the existing results can be generalized to cases where the individuals refine according to (1.3). The general answer is negative: in a finite population of Bayesian agents, communication according to a fair protocol, and information refinement according to (1.3) do not suffice for common knowledge of the signals, or even consensus.

EXAMPLE 1.1. Consider a population of three individuals with information partitions as in Figure 1.1. Let the convex signal function assign to any non-empty $J \in \sigma(\mathcal{J})$ the number

$$f(J) = \frac{1}{\#J} \sum_{\omega \in J} f(\{\omega\}), \quad (1.4)$$

with $f(\{\omega_1\}) = f(\{\omega_2\}) = 2$, $f(\{\omega_3\}) = f(\{\omega_4\}) = 3$, $f(\{\omega_5\}) = f(\{\omega_6\}) = 1$, and suppose that they refine their partitions according to Equation (1.3). Let the true state be ω_1 , and consider the fair protocol: 1 talks to 2, who talks to 3, who talks to 1, and so on. Before they start communicating there is no consensus as $f_1(\omega_1) = f_3(\omega_1) = 2 \neq 5/2 = f_2(\omega_1)$. When 1 says “2”, 2 does not learn anything since $R_1 = \Omega$ is $\sigma(\mathcal{I}_2)$ -measurable, and therefore does not refine \mathcal{I}_1 . Similarly, 3 does not refine \mathcal{I}_3 since $R_2 = \{\omega_1, \dots, \omega_4\} \in \sigma(\mathcal{I}_3)$, and 1 does refine \mathcal{I}_1 since $R_3 = \{\omega_1, \dots, \omega_4\} \in \sigma(\mathcal{I}_1)$. Therefore, they will never reach a consensus.

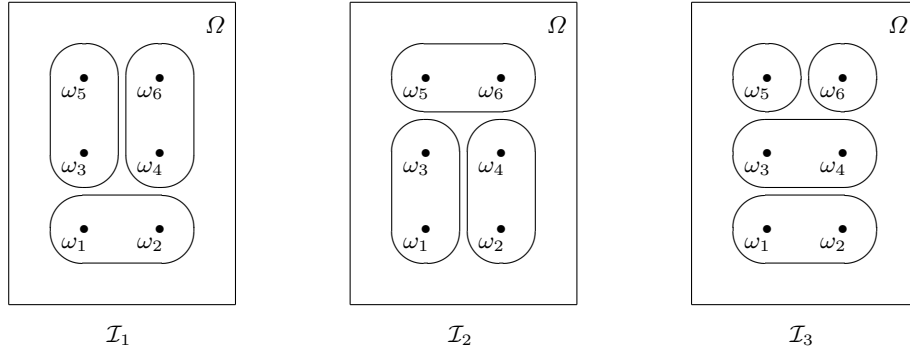


Fig. 1.1. Updating given the actual signal does not lead to consensus.

Though a general result cannot be established, we can show that there is a class of fair protocols which lead to consensus even when the individuals update their information given the actual signals, i.e., when they refine their partitions according to (1.3). The protocols that we restrict our attention to, are those which satisfy information exchange (Krasucki, 1996).

PROPOSITION 1.2. *Consider a population of like-minded individuals who refine their information given the actual signal, i.e., as in Equation (1.3). If the signal function is union consistent, and the protocol is fair and satisfies information exchange, a consensus will be reached.*

The previous proposition generalizes the one proven by Krasucki (1996), who showed that a consensus is achieved in fair protocols that satisfy information exchange whenever the individuals update their information given all the hypothetical signals, i.e., whenever they refine their partitions according to (1.2). In the next section we provide sufficient conditions for consensus when communication takes place according to an arbitrary fair protocol.

1.4 Common knowledge of the protocol and consensus

As we have already mentioned above, when communication is private, the only agent who updates is the recipient of the signal. However, when 1 talks to 2 privately, it does not mean that 3 does not know that this conversation has taken place. It might be the case

that 3 does not know what 1 has said, but knows that something has been said. That is, one has to distinguish between knowing *that* something has been said and knowing *what* this “something” is.

In this section we focus on the first case when the structure of the protocol is commonly known, i.e., when it is commonly known who talks to whom at every period. Recall the example presented in the introduction (Heifetz, 1996): when Alice whispers something in Bob’s ear, Carol is present, and therefore the fact that they – Alice and Bob – have talked is commonly known. It would not be the case, had Carol had her eyes shut. In this case, not even the fact that they had talked would have been known.

Suppose that the structure of the protocol is commonly known in the population at every period in time. In this case one can explicitly assume that the third party (Carol) can make some inference about the signal that the sender (Alice) has sent to the recipient (Bob), and then condition on the set of signals that do not contradict her observations. Recall Example 1.1: When 2 talks to 3, 1 (the third party) considers as possible the states that induce a $\sigma(\mathcal{I}_3)$ -measurable signal transmitted by 2. That is, the signal sent by 2 (the sender) must be consistent (measurable) with the information partition of 3 (the recipient) after the communication has taken place. In this specific example it is easy to see that 1 does not learn anything as all signals that could have been sent by 2 are $\sigma(\mathcal{I}_3)$ -measurable: if 2 said “5/2” then 3 would condition with respect to $\{\omega_1, \dots, \omega_4\}$ which is $\sigma(\mathcal{I}_3)$ -measurable, whilst if 2 said “1” then 3 would condition on $\{\omega_5, \omega_6\}$ which is also $\sigma(\mathcal{I}_3)$ -measurable. Therefore, any signal would have been consistent with 3’s partition, and therefore 1 the third party conditions on Ω . Thus, she does not learn anything from knowing that 2 has talked to 3.

Formally, if the protocol is commonly known third parties, i.e., individuals other than the sender and the receiver consider as possible all those signals that are consistent with the revised information partition of the recipient, i.e., all elements of $\mathcal{P}_{s_t}^t$ that are $\sigma(\mathcal{I}_{r_t}^{t+1})$ -measurable. That is, knowing the structure of the protocol, i.e., knowing that s_t talked to r_t at time t , allows third parties to condition on

$$S_{s_t}^t = \{\omega \in \Omega : P_{s_t}^t(\omega) \in \sigma(\mathcal{I}_{r_t}^{t+1})\}. \quad (1.5)$$

It is straightforward that $R_{s_t}^t \subseteq S_{s_t}^t$, i.e., third parties receive and process less information than the receiver. In addition, the true signal is never ruled out by anyone. This follows from the fact that the true signal is what has actually been said by s_t , and therefore it is necessarily measurable with respect to the recipient's partition.

Using this knowledge, every $j \in N$ refines her information partition at time t according to the rule

$$\mathcal{I}_j^{t+1} = \begin{cases} \mathcal{I}_j^t \vee \mathcal{S}_{s_t}^t & \text{if } j \neq r_t, \\ \mathcal{I}_j^t \vee \mathcal{R}_{s_t}^t & \text{if } j = r_t, \end{cases} \quad (1.6)$$

where $\mathcal{S}_{s_t}^t = \{S_{s_t}^t, (S_{s_t}^t)^c\}$. This refining mechanism implies that individuals make use of their knowledge about the structure of the protocol.

Notice that when the partitions are refined according to (1.6), the receiver still refines according to the actual signals. Third parties receive some aggregate information about the actual signal, which of course is less informative than the actual signal. This assumption reflects the idea that since they do not hear the actual signal, they do not learn as much as the receiver. However, they condition on what they learn, and not on all the possible things that they could have learned under all different hypothetical scenarios. Thus, they also update according to actual and not hypothetical information.

Note also that we require common knowledge of the protocol, in order everybody to know that third parties have refined their information due their knowledge of the protocol, and everybody to know that everybody knows that third parties have refined, and so on. Then we can actually provide a general consensus result for all fair protocols.

PROPOSITION 1.3. *Consider a population of like-minded individuals, who communicate the value of a convex function according to some fair protocol. If the protocol is commonly known, i.e., if agents refine their information according to (1.6), a commonly known consensus will be achieved.*

Remember Example 1.1, and assume updating as in (1.6). When individual 1 hears 3's signal, she does not refine her information. Individual 2, on the other hand, infers that the only signal that 1 could have heard is "2": the other signals, "1" and "3", would have led 1 to refine her partition. Exploiting this information, 2 refines her information partition.

If everybody updates according to (1.6), it is easy to verify that eventually they will agree on the commonly known signal “2”.

Appendix

PROOF OF PROPOSITION 1.2. For all $t \in \mathbb{N}$ and all $i \in N$: $\mathcal{I}_i^{t+1} \subseteq \sigma(\mathcal{J})$ is weakly finer than $\mathcal{I}_i^t \subseteq \sigma(\mathcal{J})$. As $\sigma(\mathcal{J})$ is finite, there is a $T \in \mathbb{N}$ after which no refinement occurs: for each $i \in N$ there is a partition \mathcal{I}_i^* such that $\mathcal{I}_i^t = \mathcal{I}_i^*$ for all $t > T$. Hence, for each $i \in N$ there is a set $R_i^* \subseteq \Omega$ such that $R_i^t = R_i^*$ for all $t > T$. Let $i, j \in N$ be connected. By fairness and information exchange, there are $t, t' > T$ with $s_t = r_{t'} = i$ and $r_t = s_{t'} = j$. Since no refinement occurs after T , it follows that $R_i^* \in \sigma(\mathcal{I}_j^*)$. Moreover, $R_i^* \in \sigma(\mathcal{I}_i^*)$ by definition. Hence $R_i^* \in \sigma(\mathcal{I}_i^*) \cap \sigma(\mathcal{I}_j^*) = \sigma(\mathcal{I}_i^* \wedge \mathcal{I}_j^*)$. Let $\omega \in \Omega$ be the true state. As $\omega \in R_i^*$ and $R_i^* \in \sigma(\mathcal{I}_i^* \wedge \mathcal{I}_j^*)$, non-delusion implies that the element $(I_i^* \wedge I_j^*)(\omega)$ of the partition $\mathcal{I}_i^* \wedge \mathcal{I}_j^*$ containing ω satisfies $(I_i^* \wedge I_j^*)(\omega) \subseteq R_i^*$. Similarly, $(I_i^* \wedge I_j^*)(\omega) \subseteq R_j^*$, so

$$(I_i^* \wedge I_j^*)(\omega) \subseteq R_i^* \cap R_j^*,$$

i.e., the signals of i and j satisfy pairwise common knowledge at ω . From Cave (1983) and Bacharach (1985), it follows that the signals that i and j transmit must be the same. As this holds for all connected pairs and the protocol is fair, it follows that all agents have the same signal, i.e., there is consensus. *QED.*

PROOF OF PROPOSITION 1.3. Step 0 Similarly to Proposition 1.2, there is $T > 0$ such that $\mathcal{I}_i^t = \mathcal{I}_i^*$ for every $t > T$, and every $i \in N$.

Step 1 Since the protocol is fair, there is a path of directed edges which starts from 1, passes from every individual (possibly more than once), and returns to 1. Let the finite sequence of individuals $\{p_1, \dots, p_m\} \in N^m$, with $m \geq n + 1$, determine this path:

- $p_1 = p_m = 1$,
- for every $j = 1, \dots, m - 1$ there is a directed edge from p_j to p_{j+1} which belongs to the path, i.e., there are infinitely many $t > T$ such that $s_t = p_j$ and $r_t = p_{j+1}$.

Let $t_1 > T$ be such that $s_{t_1} = p_1 = 1$ and $r_{t_1} = p_2$. Let now t_2 be the smallest $t > t_1$, such that $s_{t_2} = p_2$ and $r_{t_2} = p_3$. In general, let t_j be the smallest $t > t_{j-1}$, such that $s_{t_j} = p_j$

and $r_{t_j} = p_{j+1}$. For simplicity and without loss of generality let $m - 1 = n$, i.e., every individual appearing only once in the path. Let also for notation simplicity $p_i = i$ for all $i = 1, \dots, m - 1$. Clearly (because of Step 0), the information partition of i at t_i is given by \mathcal{I}_i^* for every $i \in N$.

Step 2 Consider an arbitrary $i \in N$. By definition, $S_i = \{\omega' \in \Omega : P_i(\omega') \in \sigma(\mathcal{I}_{i+1}^*)\}$ is $\sigma(\mathcal{I}_{i+1}^*)$ -measurable. Since the structure of the protocol is known and no individual refines after T , it follows from (1.6) that $S_i \in \sigma(\mathcal{I}_j^*)$ for every $j \in N \setminus \{i + 1\}$. Thus,

$$S_i \in \bigcap_{j \in N} \sigma(\mathcal{I}_j^*) = \sigma(\bigwedge_{j \in N} \mathcal{I}_j^*) = \sigma(\mathcal{M}^*),$$

which (given non-delusion) implies that $M^*(\omega) \subseteq S_i$. Hence, $P_i(\omega') \in \sigma(\mathcal{I}_{i+1}^*)$, for every $\omega' \in M^*(\omega)$. In addition, since $M^*(\omega) \in \sigma(\mathcal{I}_{i+1}^*)$, it follows that

$$P_i(\omega') \cap M^*(\omega) \in \sigma(\mathcal{I}_{i+1}^*),$$

for every $\omega' \in M^*(\omega)$.

Step 3 Let

$$\omega_1^0 \in \arg \min_{\omega' \in M^*(\omega)} f(P_1(\omega') \cap M^*(\omega)).$$

It follows from Step 2 that $P_1(\omega_1^0) \cap M^*(\omega) \in \sigma(\mathcal{I}_2^*)$. Hence, there is a finite collection $F(\omega_1^0) = \{\omega_2^1, \dots, \omega_2^{J_2}\} \subseteq P_1(\omega_1^0) \cap M^*(\omega)$ such that

- $I_2^*(\omega_2^j) \in \mathcal{I}_2^*$, for all $j = 1, \dots, J_2$,
- $I_2^*(\omega_2^j) \cap I_2^*(\omega_2^k)$, for all $j \neq k$, and
- $P_1(\omega_1^0) \cap M^*(\omega) = I_2^*(\omega_2^1) \cup \dots \cup I_2^*(\omega_2^{J_2})$.

Let

$$\omega_2^0 \in \arg \min_{\omega' \in F(\omega_1^0)} f(I_2^*(\omega')) = \arg \min_{\omega' \in F(\omega_1^0)} f(P_2(\omega')) = \arg \min_{\omega' \in F(\omega_1^0)} f(P_2(\omega') \cap M^*(\omega)).$$

We iteratively define ω_i^0 for every $i = 2, \dots, m$.

Step 4 We consider the following (exhaustive and mutually exclusive) cases:

- (I) for some $i \in \{2, \dots, m\}$, there are $\omega_i^j, \omega_i^k \in F(\omega_{i-1}^0)$ such that $f(I_i^*(\omega_i^j)) \neq f(I_i^*(\omega_i^k))$,
and
- (II) for all $i \in \{2, \dots, m\}$, and for all $\omega_i^j, \omega_i^k \in F(\omega_{i-1}^0)$, $f(I_i^*(\omega_i^j)) = f(I_i^*(\omega_i^k))$.

First we show that (I) cannot occur. It follows from convexity and from Step 3 that

$$\begin{aligned}
f(P_{i-1}(\omega_{i-1}^0) \cap M^*(\omega)) &= \sum_{j=1}^{J_i} \alpha_j f(I_i^*(\omega_i^j)) \\
&\geq \min_{\omega' \in F(\omega_{i-1}^0)} f(I_i^*(\omega')) \\
&= f(I_i^*(\omega_i^0)) \\
&= f(P_i(\omega_i^0) \cap M^*(\omega)),
\end{aligned}$$

with the (strict) inequality holding if and only if i is such that there are $\omega_i^j, \omega_i^k \in F(\omega_{i-1}^0)$ such that $f(I_i^*(\omega_i^j)) \neq f(I_i^*(\omega_i^k))$. However, if there is some $i = 2, \dots, m$ that satisfies (I), then

$$f(P_1(\omega_1^0)) > f(P_m(\omega_m^0)),$$

which is a contradiction, since $p_1 = p_m$, and $\omega_m^0 \in M^*(\omega)$. Hence, case (II) necessarily occurs.

In case (II), by definition, $P_i(\omega_i^j) = P_i(\omega_i^k)$ for all $\omega_i^j, \omega_i^k \in F(\omega_{i-1}^0)$. Hence, $P_i(\omega_i^j) = P_i(\omega_i^0)$ for every $\omega_i^j \in F(\omega_{i-1}^0)$, implying that

$$P_{i-1}(\omega_{i-1}^0) \cap M^*(\omega) \subseteq P_i(\omega_i^0) \cap M^*(\omega),$$

for every $i = 2, \dots, m$. It follows that $P_i(\omega_i^0) \cap M^*(\omega) = P^*$ for all $i \in N$. It follows then (from Step 2) that

$$P^* \in \bigcap_{i \in N} \sigma(\mathcal{I}_i^*) = \sigma(\bigwedge_{i \in N} \mathcal{I}_i^*) = \sigma(\mathcal{M}^*).$$

Given non-delusion, everybody's signal is commonly known. It follows then from Cave (1983) and Bacharach (1985) that the signals have converged to a commonly known consensus. *QED.*

Aggregate information, common knowledge and agreeing not to bet

2.1 Introduction

Aumann (1976) was the first one to formalize the concept of common knowledge. In his seminal paper he showed that if two people with a common prior have commonly known posterior probability assessments about an event, these probabilities are identical. Geanakoplos and Polemarchakis (1982) introduced the problem into a dynamic framework, by showing that if two individuals communicate their probability assessments back and forth and update accordingly, they will eventually agree on a common posterior probability.

Sebenius and Geanakoplos (1983) extended Aumann's result to expectations, by proving that if two people's expectations about a random variable are commonly known, then they are necessarily equal. A number of generalizations appeared in the literature ever since (Nielsen et al., 1990; Nielsen, 1995; Hanson, 1998; Hanson, 2002). A direct application of this proposition is the famous no-bet theorem, which states that two risk-averse individuals with a common prior will never agree to participate in a gamble, if their willingness to bet is commonly known. Milgrom and Stokey (1980) had already addressed this problem, by showing that common knowledge precludes trading among risk-averse agents in an uncertain environment.

This note extends the no-bet theorem to cases where the gamble can take place even if not everybody is willing to participate. Consider for instance a soccer game between A and B, and suppose that Alice predicts that A will win, Bob predicts that B will win, while Carol predicts a draw. In order to participate in the bet, they have to pay a one-

dollar entry fee and the one who predicts correctly takes it all. If however Carol refuses to participate in the gamble, while Ann and Bob accept, the bet can still take place, with the prize being equal to 2, instead of 3 dollars now. The only difference is that if a draw occurs nobody will win the prize and they will receive their entry fees back.

There is an important feature in this type of bets. The willingness to participate depends, not only on the private information, but also on who you are playing against. In the previous example, suppose that Alice believes that the probability of a draw is higher than A's victory. Then, if Carol had stayed in the bet, Alice would have rejected participation. The reason why she accepts to gamble against Bob's pick is that she believes that the probability of B winning is lower than the probability of A winning.

What is not very clear from the previous analysis is the answer to the following question: what would Ann do if she knew that one more person was willing to gamble but she did not know who? In this case, she would form beliefs given her private information about who the other player was and she would maximize her expected payoff in a Bayesian manner. However, if she knew that everybody knew how many people were willing to participate, her beliefs about who the other player was would depend on what she believed about the other two people's beliefs, and so on.

In this note I show that the bet will not take place if the number of people who are willing to participate is commonly known. This result is quite surprising, since the expected payoff, and therefore the decision about whether to participate or not depends on the identity of the other participants. However, for the bet not to take place it is sufficient to have common knowledge of the aggregate behavior (how many people participate), instead of the individual behavior (who are the other participants).

This result can be easily extended to negative-sum bets. Consider for instance a lottery where the prize depends on the number of participants, but not on who participates. In this case, if nobody has bought the winning coupon, the participants do not receive their entry fee back. Then it is straightforward that common knowledge of the number of participants precludes the bet.

2.2 Knowing how many players participate and agreeing not to bet

2.2.1 Information and knowledge

Consider a finite state space Ω , and a population $N = \{1, \dots, n\}$. The measure π determines the (common) prior beliefs of the individuals in the population over Ω and is assumed to assign positive probability to every state: $\pi(\omega) > 0$ for every $\omega \in \Omega$. Every individual is endowed with an *information partition* \mathcal{P}_i over Ω . The set $P_i(\omega) \in \mathcal{P}_i$ contains the states that i cannot distinguish from ω , with ω itself being one of those. Let $\mathcal{J} = \bigvee_{i=1}^n \mathcal{P}_i$, and $\mathcal{M} = \bigwedge_{i=1}^n \mathcal{P}_i$ denote the join (coarsest common refinement), and the meet (finest common coarsening) of the information partitions respectively. We define knowledge as usual, i.e., we say that i *knows* some $B \subseteq \Omega$ at ω whenever $P_i(\omega) \subseteq B$. The event B is *commonly known* if $M(\omega) \subseteq B$, where $M(\omega)$ denotes the member of \mathcal{M} that contains ω .

2.2.2 Gambles with limited participation

We define a *gamble* that allows for limited participation as a partition $\mathcal{G} = \{G_1, \dots, G_n\}$ of Ω , where G_i denotes the set of states where i wins. Participating in the gamble has a fixed cost (entry fee), which for simplicity and without loss of generality is normalized to 1 unit. Let $A_i = \{0, 1\}$ denote i 's action space: i plays 1 when she is willing to participate in the bet and 0 when she is not. Let the *action function* $a_i : \Omega \rightarrow A_i$ determine the action that player i undertakes at every state ω . A natural assumption is that a_i is $\sigma(\mathcal{P}_i)$ -measurable, i.e., $a_i(\omega') = a_i(\omega)$ for all $\omega' \in P_i(\omega)$, implying that i knows what she is doing. Let

$$S(\omega) = \{i \in N : a_i(\omega) = 1\} \subseteq N \quad (2.1)$$

denote the set of people who agree to participate while being at ω , and $s(\omega)$ the cardinality of $S(\omega)$. If $i \in S(\omega)$ wins the bet, she receives $s(\omega)$ units – 1 from each participant. If on the other hand $j \notin S(\omega)$ wins, no units are transferred among the participants. Formally, i 's payoff at ω , depends on $S(\omega)$, i.e., on who participates at this state, and is equal to

$$U_i(\omega) = (s(\omega)1_{\{\omega \in G_i\}} - 1_{\{\omega \in \bigcup_{j \in S(\omega)} G_j\}})1_{\{i \in S(\omega)\}}, \quad (2.2)$$

where $1_{\{\cdot\}}$ denotes the indicator function. If i does not participate in the bet, i.e., if $i \notin S(\omega)$, her payoff is equal to 0. If i on the other hand participates, i.e., if $i \in S(\omega)$, her payoff is equal to $s(\omega) - 1$ if she wins, -1 if another participant wins and 0 if the winner has chosen not to participate.

Player i 's expected payoff at ω is equal to

$$E[U_i|P_i(\omega)] = \sum_{\omega' \in P_i(\omega)} \pi(\omega'|P_i(\omega)) \left(s(\omega')\pi(G_i|\omega') - \sum_{j \in S(\omega')} \pi(G_j|\omega') \right), \quad (2.3)$$

when i participates in the bet, and 0 otherwise. Let R_i be the set of states where i 's expected payoff is strictly positive:

$$R_i = \{\omega \in \Omega : E[U_i|P_i(\omega)] > 0\}. \quad (2.4)$$

Then we say that i is **risk-averse** at ω whenever the following condition holds: $\omega \in R_i$ if and only if $i \in S(\omega)$. In other words, a risk-averse individual gambles if and only if her expected payoff is strictly positive. We assume that all individuals are risk-averse at every $\omega \in \Omega$.

2.2.3 Main result

It is commonly known at ω that s individuals participate, if $s(\omega') = s$ for every $\omega' \in M(\omega)$. Clearly, common knowledge of how many people participate is weaker than common knowledge of who participates. If s is commonly known, i does not necessarily know who she is playing against at ω . If on the other hand S is commonly known, then $S(\omega') = S$ for all $\omega' \in M(\omega)$, and it is straightforward that no bet can take place (it follows from Sebenius and Geanakoplos, 1983). The following result states that it suffices to require common knowledge of s – instead of S – for the bet not to take place.

THEOREM 2.1. *Consider a gamble that allows for limited participation and suppose that all individuals are risk-averse at every state. Then no bet can take place, if it is commonly known how many agents are willing to accept the bet.*

This result is quite surprising since the payoff at ω depends on $S(\omega)$, rather than $s(\omega)$. Therefore, it is not obvious that common knowledge of s precludes the bet from taking place.

It is straightforward to extend this result to negative sum bets, where the payoff depends on the number of participants but not on who participates:

$$\hat{U}_i(\omega) = (s(\omega)1_{\{\omega \in G_i\}} - 1)1_{\{i \in S(\omega)\}}.$$

In this case, if i participates, she has to pay the 1 unit entrance fee, even if the winner decides to stay out. Clearly, $\hat{U}_i(\omega) \leq U_i(\omega)$ for all $\omega \in \Omega$, and therefore i participates less often than she would do in a zero-sum bet, as the one analyzed above. It can be shown then that a bet that allows for limited participation cannot take place if s is commonly known. Formally, the proof is identical to Step 2 of the proof of Theorem 2.1.

Appendix

PROOF OF THEOREM 2.1. Step 1 It follows from common knowledge that $s(\omega) = s$ for every $\omega \in M(\omega_0)$, where ω_0 denotes the actual state. It follows from risk-aversion that i participates at $\omega \in M(\omega_0)$ if and only if

$$E[U_i|P_i(\omega)] = s\pi(G_i|P_i(\omega)) - \sum_{\omega' \in P_i(\omega)} \pi(\omega'|P_i(\omega)) \sum_{j \in S(\omega')} \pi(G_j|\omega') > 0. \quad (2.5)$$

Let W be the set of states where the winner is willing to participate, i.e., $W \subseteq \Omega$ is the set of states that satisfy the following condition: if $\omega \in G_i$, then $i \in S(\omega)$. Then, we rewrite (2.5) as follows:

$$E[U_i|P_i(\omega)] = E[U_i|P_i(\omega) \cap W]\pi(P_i(\omega) \cap W) + E[U_i|P_i(\omega) \cap W^c]\pi(P_i(\omega) \cap W^c) > 0. \quad (2.6)$$

Notice that everybody's expected payoff is equal to 0 when the winner does not participate. Hence, $E[U_i|P_i(\omega) \cap W^c] = 0$, implying that

$$E[U_i|P_i(\omega)] = E[U_i|P_i(\omega) \cap W]\pi(P_i(\omega) \cap W) > 0. \quad (2.7)$$

Thus i participates at ω if and only if $\pi(P_i(\omega) \cap W) > 0$ and $E[U_i|P_i(\omega) \cap W] > 0$, which given Equation (2.5), is equivalent to

$$s\pi(G_i|P_i(\omega) \cap W) - \sum_{\omega' \in P_i(\omega) \cap W} \pi(\omega'|P_i(\omega) \cap W) \sum_{j \in S(\omega')} \pi(G_j|\omega') > 0. \quad (2.8)$$

Since, $\omega' \in W$ it follows that there is one $k \in S(\omega')$ such that $\omega' \in G_k$. Hence, $\sum_{j \in S(\omega')} \pi(G_j | \omega') = 1$ for every $\omega' \in P_i(\omega) \cap W$. Therefore, for every $\omega \in M(\omega_0)$, i participates if and only if

$$s\pi(G_i | P_i(\omega) \cap W) > 1. \quad (2.9)$$

Step 2 If $s = 1$, the proof is straightforward, since $U_i(\omega) \leq 0$ for every $i \in N$ and for all $\omega \in M(\omega_0)$ and therefore no i participates.

Suppose now that $s \geq 2$. It follows from (2.7) and (2.9) that

$$s\pi(G_i | P_i(\omega) \cap W)\pi(P_i(\omega) \cap W | M(\omega_0) \cap W) > \pi(P_i(\omega) \cap W | M(\omega_0) \cap W), \quad (2.10)$$

for all $\omega \in R_i \cap M(\omega_0)$. Obviously, $R_i \cap M(\omega_0)$ is $\sigma(\mathcal{P}_i)$ -measurable. It follows then from summing over $R_i \cap M(\omega_0)$ that

$$s \sum_{P_i \subseteq R_i \cap M(\omega_0)} \pi(G_i | P_i \cap W)\pi(P_i \cap W | M(\omega_0) \cap W) > \sum_{P_i \subseteq R_i \cap M(\omega_0)} \pi(P_i \cap W | M(\omega_0) \cap W). \quad (2.11)$$

Since $R_i \cap M(\omega_0) \subseteq M(\omega_0)$ it follows that

$$s \sum_{P_i \subseteq M(\omega_0)} \pi(G_i | P_i \cap W)\pi(P_i \cap W | M(\omega_0) \cap W) \geq s \sum_{P_i \subseteq M(\omega_0)} \pi(P_i \cap W | M(\omega_0) \cap W). \quad (2.12)$$

It follows then from (2.11) and (2.12) that

$$s\pi(G_i | M(\omega_0) \cap W) > \pi(R_i \cap W | M(\omega_0) \cap W). \quad (2.13)$$

It is straightforward that

$$\pi(R_i \cap W | M(\omega_0) \cap W) = \pi(R_i | M(\omega_0) \cap W) = \sum_{\omega \in R_i} \pi(\omega | M(\omega_0) \cap W). \quad (2.14)$$

Then it follows from (2.13) and (2.14) that

$$s\pi(G_i | M(\omega_0) \cap W) > \sum_{\omega \in R_i} \pi(\omega | M(\omega_0) \cap W). \quad (2.15)$$

Summing over the individuals entails

$$s \sum_{i \in N} \pi(G_i | M(\omega_0) \cap W) > \sum_{i \in N} \sum_{\omega \in R_i} \pi(\omega | M(\omega_0) \cap W), \quad (2.16)$$

which in turn yields the contradiction $s > s$, since (i) all winners participate in $M(\omega_0) \cap W$ implying that $\sum_{i \in N} \pi(G_i | M(\omega_0) \cap W) = 1$, and (ii) at every $\omega \in M(\omega_0) \cap W \subseteq M(\omega_0)$ exactly s individuals participate. Hence, the probability $\pi(\omega | M(\omega_0) \cap W)$ appears s times in the sum. In addition, all states $\omega \in M(\omega_0) \cap W \subseteq M(\omega_0)$ appear in the sum, implying that it is equal to s , which completes the proof. *QED.*

Testing rationality on primitive knowledge

3.1 Introduction

One of the central components of modern economics is the consideration that agents are limited in their means to access and treat information. The “bounded rationality” literature, which specifically addresses this issue, dates back to Simon (1955), and has surged in the last decades (Aumann, 1997; Lipman, 1995). Rationality and its lack thereof can be observed both in the agents’ decisions, and their information processing. Failures of rationality in the information processing are more fundamental than in decision taking, since non-rational information processing entails non-rational decision making, whilst the converse is not necessarily true.

Hintikka (1962), Aumann (1976) and Geanakoplos (1989) introduced a semantic model of information structures that represent the information processing of both perfectly, and non-perfectly rational agents respectively (see also Brandenburger et al., 1992; and Dekel et al., 1998). It is commonly argued that the possibility correspondence of a perfectly rational agent has to be partitional, since the agent can exclude all the states with different information, and no others. Hence, rational agents should exhibit partitional possibility correspondences, and non-partitional possibility correspondences should be taken as a sign of irrationality (see, e.g, chapter 3 of Rubinstein, 1998).

An important question is whether one can test whether an agent’s possibility correspondences is partitional or not. Recall that, by definition, an agent’s information is partitional if, when the agent being in state ω considers the state ω' as possible, the same agent being at ω' would also consider ω as possible. Testing this is problematic, since it

requires to observe the agent's knowledge in both states ω and ω' . Hence, the semantic model provides no direct way of testing whether an agent processes information rationally, or not.

Kripke (1963), Bacharach (1985) and Samet (1990) introduced a syntactic model of knowledge in the form of a system of propositional calculus. A syntactic model explicitly describes the set of propositions known by the agent at each state.

A syntactic model of knowledge defines a semantic model in a natural way: the states the agent believes as possible are the states in which all the known propositions are true. A semantic model also defines a syntactic model: the known propositions are supersets of the set of states considered as possible. Because they allow for a variety of propositions and more elaborate state spaces, syntactic models are a more general framework than semantic ones.

In syntactic models, Bacharach (1985), and Samet (1990) independently showed that whenever the agent's knowledge satisfies the basic axioms of 1) knowledge, 2) positive introspection, and 3) negative introspection, the induced possibility correspondence is partitional. The axiom of knowledge says that if the agent knows some proposition, then this proposition is true. Under positive introspection, if the agent knows a proposition, then she also knows that she knows this proposition. Negative introspection says that if the agent does not know a proposition, she knows that she does not know it. Both the knowledge and positive introspection axioms are based on positive knowledge. On the other hand, negative introspection is based on knowledge of oneself's ignorance.

We adopt the commonly accepted view that the knowledge and positive introspection axioms are rather non-problematic, and can be assumed if necessary. On the other hand, the negative introspection axiom is more controversial (Geanakoplos, 1989; Lipman, 1995), and should be accepted or rejected based on a combination of empirical evidence and logical implications.

Unlike partitional possibility correspondences in the semantic model, the axioms of knowledge, positive introspection and negative introspection are defined state by state. Therefore, observation of the agent's knowledge in particular states is enough to test whether these axioms are satisfied are not. Hence, syntactic models provide an appropriate framework in which the agent's information processing abilities can be tested.

The main difficulty arising from testing negative introspection is the infinite cardinality of the set of propositions. Indeed, every primitive proposition (fact) generates a sequence of epistemic (modal) propositions. We introduce an axiom of positive negation under which, if a proposition is known, it is known that its negation is not known. Positive negation, like knowledge and positive introspection, is based on positive knowledge by the agent, and is defined state by state. Our main result, Theorem 3.5, shows that under knowledge, positive introspection, and positive negation, negative introspection holds if and only if it holds for the primitive propositions. Hence, negative introspection is testable, and is sufficient for partitional information structures.

It is particularly interesting to look at the implications of our theorem when knowledge in the syntactic model stems from a semantic model, a condition that we call semantic knowledge. This is the natural setup to consider when the semantic model forms a complete description of the agent's knowledge: propositions known in the syntactic model are the ones arising from the semantic model.

We show in Proposition 3.9 that positive negation is implied by semantic knowledge. Furthermore, in this framework, partitional information and negative introspection are known to be equivalent. It follows that testing partitional information structures is equivalent to testing negative introspection, which is in turn equivalent to testing negative introspection for the primitive propositions.

To summarize our results, observation of primitive knowledge is enough to test negative introspection under positive negation, and is enough to test partitional information structures when knowledge comes from a semantic model. We find the latter particularly striking, as partitional information structures is a property defined on the information structure as a whole and for the whole set of propositions, and yet, it can be tested 1) separately on each state and 2) through knowledge of primitive propositions only.

3.2 Knowledge

We recall the standard model of propositional calculus to model knowledge (Kripke, 1963; Bacharach, 1985; Samet, 1990). Let Φ be the countable set of *propositions*, which de-

scribe the relevant characteristics of the environment. The mappings $\kappa : \Phi \rightarrow \Phi$ and $\neg : \Phi \rightarrow \Phi$ stand for the propositions “the agent **knows** ϕ ”, and “**not** ϕ ” respectively.

Let Ω_0 denote the set of all mappings $\omega : \Phi \rightarrow \{0, 1\}$ that satisfy complementarity: $\omega(\phi) = 1$ if and only if $\omega(\neg\phi) = 0$, for every $\phi \in \Phi$. We say that a proposition ϕ is true at a state $\omega \in \Omega_0$, and we write $\phi \in \omega$, if and only if $\omega(\phi) = 1$. Alternatively, a state $\omega \in \Omega_0$ can be identified by the set $\{\phi \in \Phi : \omega(\phi) = 1\}$, i.e., by the propositions that are true in this state. The ken (set of known propositions) at a state ω is defined as follows:

$$K(\omega) = \{\phi \in \Phi : \kappa\phi \in \omega\}. \quad (3.1)$$

Now consider a subset $\Omega \subset \Omega_0$. A state $\omega' \in \Omega$ is considered as possible while being at $\omega \in \Omega$ if every known proposition at ω is true at ω' . That is, the **possibility correspondence** $P : \Omega \rightarrow 2^\Omega$, maps every state ω to the set of states considered as possible by the agent while being at ω :

$$P(\omega) = \{\omega' \in \Omega : K(\omega) \subseteq \omega'\}. \quad (3.2)$$

A possibility correspondence is called **partitional** whenever $P(\omega) = P(\omega')$ for every $\omega' \in P(\omega)$, and every $\omega \in \Omega$.

The three fundamental axioms of propositional calculus are:

- (K_1) if $\kappa\phi \in \omega$ then $\phi \in \omega$ (axiom of knowledge),
- (K_2) if $\kappa\phi \in \omega$ then $\kappa\kappa\phi \in \omega$ (positive introspection),
- (K_3) if $\neg\kappa\phi \in \omega$ then $\kappa\neg\kappa\phi \in \omega$ (negative introspection).

Samet (1990) defines the following state spaces:

- Ω_1 : the set of states that satisfy (K_1),
- Ω_2 : the set of states that satisfy (K_1), (K_2),
- Ω_3 : the set of states that satisfy (K_1), (K_2), and (K_3),

and proves the following result:

PROPOSITION 3.1 (SAMET, 1990). *If $\Omega \subseteq \Omega_3$, the possibility correspondence P is partitional.*

The previous proposition follows from the fact that for $\Omega \subseteq \Omega_3$, $\omega' \in P(\omega)$ if and only if $K(\omega) = K(\omega')$, also proven by Samet (1990). In other words, $(K_1) - (K_3)$ imply that the possibility correspondence is such that the states considered as possible at ω are those that yield exactly the same knowledge as ω .

3.3 Primitive propositions and negative introspection

Proposition 3.1 provides a way to test for partitional information structures, since $(K_1) - (K_3)$ can be tested state by state. Still, if one wishes to do this, it is necessary to check (K_3) for every $\phi \in \Phi$ at every state ω . This would be practically impossible due to the infinite cardinality of Φ . Let $\Phi_0 \subset \Phi$ denote the (finite) non-empty set of primitive propositions, which are not derived from some other proposition with the use of the knowledge operator κ . These propositions, which are also called atomic or non-epistemic or non-modal, refer to natural events (facts). Aumann (1999) defines the primitive propositions as substantive happenings that are not described in terms of people knowing something. Bacharach (1985), Samet (1990), Modica and Rustichini (1999), Hart et al. (1996) and Halpern (2001) also discuss the distinction between primitive and epistemic propositions.

For some $\phi \in \Phi$, let $B_0(\phi) = \{\phi, \neg\phi\}$, and define inductively:

$$B_n(\phi) = \{\kappa\phi', \neg\kappa\phi' \mid \phi' \in B_{n-1}(\phi)\}. \quad (3.3)$$

We call $B(\phi) = \bigcup_{n \geq 0} B_n(\phi)$ the set of propositions generated by ϕ . Obviously, the set of all propositions is given by the union of all primitive and all epistemic propositions, i.e., $\Phi = \bigcup_{\phi \in \Phi_0} B(\phi)$. The cardinality of Φ is (at least countably) infinite, implying that testing whether (K_3) holds by observing every unknown proposition would be practically impossible. We propose an alternative way to figure out whether negative introspection holds or not, by only looking at the primitive propositions. Such a task is easier, not only in terms of the cardinality of the propositions to be tested, but also in terms of complexity, since articulating high order epistemic propositions can be rather complicated.

The most disputed among the three basic axioms is negative introspection. Geanakoplos (1989) and Lipman (1995) note that negative introspection is a less realistic assumption than the other axioms (K_1) and (K_2) , since negative introspection requires the agent

to notice events that have not occurred. In general, whilst (K_1) and (K_2) are based on reasoning through knowledge, (K_3) assumes that the agent makes inference based on lack of knowledge. We introduce the following axiom:

(K_4) if $\kappa\phi \in \omega$ then $\kappa\neg\kappa\neg\phi \in \omega$ (positive negation),

which implies that if a proposition is known, then it is also known that its negation is not known. The axiom (K_4) relies on inference that can be drawn by the agent using positive knowledge (as opposed to (K_3)). We compare (K_4) with the following mild axiom of inference (Modica and Rustichini, 1994):

(K_I) if $\phi \in \omega$ implies $\phi' \in \omega$, then $\kappa\phi \in \omega$ also implies $\kappa\phi' \in \omega$ (axiom of inference).

PROPOSITION 3.2. *If $\Omega \subseteq \Omega_2$ satisfies (K_I) , then $\Omega \subseteq \Omega_4$.*

PROOF. We first show that $\kappa\phi$ implies $\neg\kappa\neg\phi$: If $\kappa\phi$ and $\kappa\neg\phi$ simultaneously hold, then (K_1) implies both ϕ and $\neg\phi$, which is a contradiction because of complementarity. Assume $\kappa\phi$. We have $\kappa\kappa\phi$ by (K_2) and $\neg\kappa\neg\phi$ by complementarity, so (K_I) yields $\kappa\neg\kappa\neg\phi$, which is the desired conclusion. *QED.*

That is, if the agent is able to make simple deductions of the form of (K_I) , then her knowledge satisfies (K_4) . We consider (K_4) to be a mild requirement on the agent's knowledge.

We define the state space

Ω_4 : the set of states that satisfy (K_1) , (K_2) , and (K_4) .

PROPOSITION 3.3. $\Omega_3 \subseteq \Omega_4$.

PROOF. Given complementarity and (K_1) , if $\kappa\phi \in \omega$, then $\neg\kappa\neg\phi \in \omega$, and therefore (K_4) follows directly from (K_3) . *QED.*

Notice that the converse does not hold, i.e., (K_4) does not imply (K_3) . Thus, $\Omega \subseteq \Omega_4$ isn't sufficient for the agent's knowledge to be partitional, as shown by the following example:

EXAMPLE 3.4. Consider the state of complete ignorance ω_0 at which nothing is known: $\neg\kappa\phi \in \omega_0$, for every ϕ . Observe that $\omega_0 \in \Omega_4$ since (K_1) , (K_2) and (K_4) are based on positive knowledge (they require some implications whenever $\kappa\phi$ is true). Now consider any state $\omega_1 \in \Omega_4$ at which it is known that some ϕ is known ($\kappa\kappa\phi$ for some ϕ). Let Ω be any state space that contains ω_0 and ω_1 . Since nothing is known at ω_0 , $K(\omega_0) = \emptyset$, and $P(\omega_0) = \Omega$. On the other hand, $\kappa\phi \in K(\omega_1)$, so that $\omega_0 \notin P(\omega_1)$. Therefore, P is not partitional.

That is, some additional condition is required to ensure that the agent's knowledge is partitional. Assuming $\Omega \subseteq \Omega_4$, our main result below offers a practical way to test for partitional information, by observing the (finitely many, and easy to articulate) primitive propositions.

THEOREM 3.5. *Consider $\Omega \subseteq \Omega_4$. Then negative introspection holds at every state for every proposition if and only if it holds for all primitive propositions.*

PROOF. It follows from Lemma 3.11 (see in the Appendix) that all states with $\kappa\phi \in \omega$ contain all $\phi' \in B(\phi)$ with an even number of negations, and all states with $\kappa\neg\phi \in \omega'$ contain all $\phi' \in B(\phi)$ with an odd number of negations. In order for negative introspection to be violated, the state must contain propositions with both even and odd number of negations. When $\kappa\neg\kappa\phi \in \omega$ and $\kappa\neg\kappa\neg\phi \in \omega$, then we consider $\phi' = \neg\kappa\phi$, and $\phi'' = \neg\kappa\neg\phi$ separately, and from the previous argument (K_3) holds at ω , which proves the theorem. *QED.*

Theorem 3.5 shows that it suffices to look at whether negative introspection holds for the primitive, rather than for all, propositions. Using this theorem, and the Bacharach-Samet result, we can corroborate that P is partitional just by fulfilling a much less demanding set of conditions than required by Proposition 3.1:

COROLLARY 3.6. *Consider $\Omega \subseteq \Omega_4$, if negative introspection holds for every primitive proposition, then P is partitional.*

3.4 Primitive propositions in semantic models of knowledge

While in syntactic models knowledge of each proposition is embodied in every state, in semantic models knowledge stems from a possibility correspondence $P : \Omega \rightarrow 2^\Omega$. An event $E \subseteq \Omega$ is known at ω , and we write $\omega \in KE$, whenever $P(\omega) \subseteq E$ (Hintikka, 1962; Aumann, 1976; Geanakoplos, 1989). That is, an event is known whenever it occurs at every contingency considered as possible. Starting from a syntactic model, and for any proposition $\phi \in \Phi$ consider the event that ϕ is true:

$$E_\phi = \{\omega \in \Omega : \phi \in \omega\}. \quad (3.4)$$

The following result relates knowledge in syntactic and semantic models.

PROPOSITION 3.7. *Consider $\Omega \subseteq \Omega_0$. If $\kappa\phi \in \omega$, then $\omega \in KE_\phi$.*

PROOF. Let $\phi \in K(\omega)$. It follows from $K(\omega) \subseteq \omega'$ that $\phi \in \omega'$, for every $\omega' \in P(\omega)$. Hence $P(\omega) \subseteq E_\phi$. *QED.*

Note however that the converse of Proposition 3.7 does not hold in general. Consider for instance the following example:

EXAMPLE 3.8. Consider the state space $\Omega \subseteq \Omega_4$:

$$\begin{aligned} \omega_1 &= \{\phi, \kappa\phi, \neg\kappa\neg\phi, \kappa\kappa\phi, \kappa\neg\kappa\neg\phi, \dots\}, \\ \omega_2 &= \{\neg\phi, \neg\kappa\phi, \kappa\neg\phi, \kappa\neg\kappa\phi, \kappa\kappa\neg\phi, \dots\}, \\ \omega_3 &= \{\phi, \neg\kappa\phi, \neg\kappa\neg\phi, \neg\kappa\neg\kappa\phi, \neg\kappa\neg\kappa\neg\phi, \kappa\neg\kappa\neg\kappa\phi, \kappa\neg\kappa\neg\kappa\neg\phi, \dots\}, \\ \omega_4 &= \{\neg\phi, \neg\kappa\phi, \neg\kappa\neg\phi, \neg\kappa\neg\kappa\phi, \neg\kappa\neg\kappa\neg\phi, \kappa\neg\kappa\neg\kappa\phi, \kappa\neg\kappa\neg\kappa\neg\phi, \dots\}. \end{aligned}$$

We have $P(\omega_1) = \{\omega_1\}$, $P(\omega_2) = \{\omega_2\}$, and $P(\omega_3) = P(\omega_4) = \{\omega_3, \omega_4\}$, so knowledge is partitional. However, the proposition $\neg\kappa\phi$ is not known at ω_3 according to the syntactic definition of knowledge ($\neg\kappa\neg\kappa\phi \in \omega_3$), although it is known according to the semantic one ($P(\omega_3) \subseteq E_{\neg\kappa\phi} = \{\omega_2, \omega_3, \omega_4\}$). In other words, although it is not known that ϕ is not known in the syntactic model, it is known that ϕ is not known in the corresponding semantic model.

Note that in the previous example, (K_3) is violated for both ϕ and $\neg\phi$ at ω_3 and ω_4 : the agent does not know that she does not know ϕ , or $\neg\phi$. The possibility correspondence is partitional despite the fact that negative introspection does not hold. This is relevant to us, since our ultimate goal is to test for partitional possibility correspondences, and testing for negative introspection is a potential way to achieve this goal.

If the converse to Proposition 3.7 is satisfied, then knowledge in the syntactic and semantic models coincide. We introduce the following axiom of semantic knowledge:

(K_E) if $\omega \in KE_\phi$, then $\kappa\phi \in \omega$ (semantic knowledge).

Under semantic knowledge, the agent's possibility correspondence is partitional if and only if negative introspection holds at every state (see Theorem 5.14 p. 177 in Chellas, 1980; or Battigalli and Bonanno, 1999). In this case, testing for negative introspection is therefore equivalent to testing for partitional possibility correspondences.

Furthermore, under semantic knowledge, positive negation is always satisfied:

PROPOSITION 3.9. *Consider $\Omega \subseteq \Omega_2$, and let (K_E) hold in Ω . Then $\Omega \subseteq \Omega_4$.*

PROOF. Consider $\kappa\phi \in \omega$. Then it follows from Proposition 3.7 that $\omega \in KE_\phi$, implying that $P(\omega) \subseteq E_\phi$. It follows from Samet (1990) that $P(\omega') \subseteq P(\omega)$, for every $\omega' \in P(\omega)$, implying that $\omega' \in KE_\phi$, for every $\omega' \in P(\omega)$. Thus $\omega' \in \backslash K \backslash E_\phi$, for every $\omega' \in P(\omega)$. Hence $\omega \in K \backslash K \backslash E_\phi$. Finally it follows from (K_E) that $K \backslash K \backslash E_\phi = E_{\kappa\neg\kappa\neg\phi}$, which concludes the proof. *QED.*

It follows from Theorem 3.5 and Proposition 3.9 that, under semantic knowledge, the agent's possibility correspondence is partitional if and only if negative introspection holds for primitive propositions.

COROLLARY 3.10. *Consider $\Omega \subseteq \Omega_2$, and assume semantic knowledge. Then P is partitional if and only if negative introspection holds for every primitive proposition.*

PROOF. We know that P is partitional if and only if: $(A_1) KE_\phi \subseteq E_\phi$, $(A_2) KE_\phi \subseteq KKE_\phi$, and $(A_3) \backslash KE_\phi \subseteq K \backslash KE_\phi$, for every $\phi \in \Phi$. It follows from $KE_\phi = E_{\kappa\phi}$ and $\backslash KE_\phi = E_{\neg\kappa\phi}$ that (A_i) is equivalent to (K_i) for every $i = 1, 2, 3$. Then the proof follows directly from Proposition 3.9. *QED.*

Starting from a semantic model, the axiom of semantic knowledge is a natural assumption, as it only requires the agent to form in the syntactic model the same inferences as those arising from the semantic model. However, Proposition 3.10 shows that semantic knowledge has a striking implication.

Indeed, partitional knowledge is by nature a semantic property. As such, it requires observation of the entire state space in order to be verified or rejected. However, when knowledge arises from a semantic model, partitional knowledge can be broken down to a state by state property, which requires conditions on the agent's knowledge on primitive propositions only. We conclude that, in these models, the agent's rational information processing is an easily testable assumption.

Appendix

LEMMA 3.11. *Consider $\Omega \subseteq \Omega_4$, and let $\phi \in K(\omega)$. Then $\phi' \in B(\phi)$ belongs to ω if and only if ϕ' contains an even number of negations.*

PROOF. Every proposition in $B_n(\phi)$ contains n knowledge operators κ and one proposition ϕ . Then, every $\phi' \in B_n(\phi)$ can be fully identified by a finite sequence of $n + 1$ variables, where i -th variable takes value 1 if there is a negation in front of the i -th knowledge operator, and 0 otherwise.

Let $B_n^e(\phi)$ denote the subset of $B_n(\phi)$ that contains an even number of negations and $B_n^o(\phi) := B_n(\phi) \setminus B_n^e(\phi)$ subset of $B_n(\phi)$ that contains an odd number of negations. There is a bijection between $B_n^e(\phi)$ and $B_n^o(\phi)$, implying that their cardinality is the same. This follows from the fact that for every $\phi' \in B_n^e(\phi)$ a unique $\phi'' \in B_n^o(\phi)$ is obtained by simply changing the $(n + 1)$ -th coordinate and vice versa.

It follows then by induction that if $\Omega \subseteq \Omega_4$ and $\kappa\phi \in \omega$, then $B_n^e(\phi) \subseteq \omega$, which implies that $B_n^e(\phi) = \omega$, thus completing the proof. *QED.*

Evolutionary Game Theory

The target projection dynamic

4.1 Introduction

The most well-known and extensively used solution concept in noncooperative game theory is the Nash equilibrium. The question how players may reach such equilibria is studied in a branch of game theory employing dynamic models of learning and strategic adjustment. The main dynamic processes in the theory of strategic form games include the replicator dynamic (Taylor and Jonker, 1978), the best-response dynamic (Gilboa and Matsui, 1991), and the Brown-Nash-von Neumann (BNN) dynamic (Brown and von Neumann, 1950). Sandholm (2005) introduced a definition for well-behaved evolutionary dynamics through a number of desiderata (see Theorem 4.6 for precise definitions):

EXISTENCE, UNIQUENESS, AND CONTINUITY OF SOLUTIONS to the specified dynamic process,

NASH STATIONARITY: the stationary points of the process coincide with the game's Nash equilibria,

POSITIVE CORRELATION: roughly speaking, the probability of “good” strategies increases, that of “bad” strategies decreases.

He showed that – unlike the replicator and the best-response dynamics – the family of BNN or excess-payoff dynamics is well-behaved.

In the present paper we analyze the *target projection dynamic* that was mentioned only briefly in the same paper by Sandholm (2005, pp. 166–167). Our main results include the following:

Although the dynamic has a certain geometric appeal, Sandholm (2005, p. 167) wrote: “Unfortunately, we do not know of an appealing way of deriving this dynamic from a model of individual choice”. This is remedied in Proposition 4.5, which provides a microeconomic foundation for the target projection dynamic. Following the control cost approach (Van Damme, 1991; Mattsson and Weibull, 2002; Voorneveld, 2006), we show that it models rational behavior in a setting where the players have to exert some effort/incur costs to deviate from incumbent strategies. In other words: the target projection dynamic is a best-response dynamic under a certain status-quo bias.

The fact that the players face control costs makes their adjustment process slower. This makes sense, since they are averse – to some extent – to deviations from their current behavior. However, despite the fact that their learning mechanism is quite conservative, the target projection dynamic is well-behaved in the sense described above. It also satisfies an additional property:

INNOVATION: if some population has not yet reached a stationary state and has unused best responses, part of the population switches to it.

This is established in Theorem 4.6. These properties imply (Hofbauer and Sandholm, 2007) that there are games where strictly dominated strategies survive under the target projection dynamic. Nevertheless, we show that strictly dominated strategies are wiped out if the “gap” between the dominated and dominant strategy is sufficiently large (Proposition 4.7) or if there are only two pure strategies (Proposition 4.8).

Like most other dynamics, the target projection dynamic belongs to family of uncoupled dynamics, where the behavior of one player is independent of payoffs to other players. Therefore, the process cannot converge to Nash equilibrium in all games (Hart and Mas-Colell, 2003). Nevertheless, some special cases can be established:

- sufficiently close to interior Nash equilibria of zero-sum games, the (standard Euclidean) distance to such an equilibrium remains constant (Corollary 4.12),
- strict Nash equilibria are asymptotically stable (Proposition 4.13),
- as are evolutionarily stable strategies (Maynard Smith, 1982) if they are interior (Corollary 4.14) or the game is 2×2 (Proposition 4.15).

The first two points rely on the analysis of stable games in Hofbauer and Sandholm (2008). The fact that in certain cases the target projection dynamic coincides with the projection dynamic, which was introduced to games and extensively studied by Lahkar and Sandholm (2008) and Sandholm et al. (2008); see Proposition 4.9.

4.2 Notation and preliminaries

4.2.1 Learning in normal form games

Consider a finite normal form game $G = (N, (A_i)_{i \in N}, (u_i)_{i \in N})$ defined as usual: $N = \{1, \dots, n\}$ is the set of players, $A_i = \{a_i^1, \dots, a_i^{J_i}\}$ is player i 's finite set of actions (pure strategies) with a_i being the typical element of A_i , and $u_i : A \rightarrow \mathbb{R}$ is i 's payoff function, where $A = \times_{i \in N} A_i$ is the game's action space. Consider i 's set of mixed strategies $\Delta_i := \{\alpha_i \in \mathbb{R}_+^{J_i} : \sum_{j=1}^{J_i} \alpha_i^j = 1\}$, with typical element $\alpha_i = (\alpha_i^1, \dots, \alpha_i^{J_i})$, and let α_i^j denote the probability that α_i assigns to a_i^j . We say that α_i is completely mixed if $\alpha_i^j > 0$ for all $j = 1, \dots, J_i$. With slight abuse of notation we write $\alpha = (\alpha_i, \alpha_{-i})$, where $\alpha_{-i} = (\alpha_1, \dots, \alpha_{i-1}, \alpha_{i+1}, \dots, \alpha_n)$. We define the expected payoff $u_i : \Delta \rightarrow \mathbb{R}$ as usual, where $\Delta := \times_{i \in N} \Delta_i$. Then, we rewrite i 's expected payoff as follows: $u_i(\alpha) = \langle \alpha_i, U_i(\alpha) \rangle$, where $\langle x, y \rangle := \sum_{i=1}^m x_i y_i$ denotes the usual inner product of two vectors $x, y \in \mathbb{R}^m$, and $U_i(\alpha) = (U_i^1(\alpha), \dots, U_i^{J_i}(\alpha))$ is the vector of expected payoffs to player i 's pure strategies given that everybody else plays according to α , i.e., $U_i^j(\alpha) = u_i(a_i^j, \alpha_{-i})$.

We say that α_i is a **best response** to α if $\langle \alpha_i, U_i(\alpha) \rangle \geq \langle \beta_i, U_i(\alpha) \rangle$ for all $\beta_i \in \Delta_i$. If α_i is a best response to α for every $i \in N$ then α is a **Nash equilibrium**. If the inequality is strict for every $i \in N$ then α is called **strict Nash equilibrium**. Obviously, strict equilibria appear only in pure strategies.

We say that an a_i^j **(weakly) dominates** a_i^k , if $U_i^j(\alpha) \geq U_i^k(\alpha)$ for all $\alpha \in \Delta$ and it holds with strict inequality for at least one α . If the inequality is strict for every $\alpha \in \Delta$ we say that a_i^j **strictly dominates** a_i^k .

Consider the standard framework of learning in normal form games (Börger and Sarin, 1997; Fudenberg and Levine, 1998; Hopkins, 2002). There is a population of individuals. Every individual repeatedly participates in G against other individuals, always holding the same role, i.e., always being the same $i \in N$. The way the individuals are matched

can vary according to the model: they play against the same opponents at every period (single-matching) or they are randomly matched with individuals playing the other roles (random-matching). At all periods every player i chooses a mixed strategy $\alpha_i \in \Delta_i$. After having observed the own payoff and (usually) the strategy profile $\alpha \in \Delta$, every $i \in N$ adjusts the own strategy according to some updating rule $\dot{\alpha}_i$, which typically puts higher weight on more profitable strategies and less weight on the ones that entail lower payoff. The adjustment rule describes the mechanism through which players learn how to play in the future, given their current observations, i.e., $\dot{\alpha} = f(\alpha)$. Dynamic processes that do not depend on the payoff earned by the other players are called uncoupled dynamics (Hart and Mas-Colell, 2003).

4.2.2 Projections

This subsection contains some results on projections. In particular, Proposition 4.1 establishes a link between maximizing a linear function and certain projection problems. Proposition 4.2 gives a simple expression for projection onto the unit simplex. The proofs are in the Appendix.

As we have already mentioned, let $\langle x, y \rangle = \sum_{i=1}^m x_i y_i$ denote the usual inner product of two vectors $x, y \in \mathbb{R}^m$. Let $\|\cdot\|$ denote the standard Euclidean norm, i.e., $\|x\| = \langle x, x \rangle^{1/2}$. For $z \in \mathbb{R}$, let $[z]_+ := \max\{z, 0\}$.

PROPOSITION 4.1. *Let $C \subseteq \mathbb{R}^n$ be nonempty and convex, $a \in \mathbb{R}^n$, $c \in C$ and $k > 0$.*

(i) *The following two claims are equivalent:*

- (a) *c solves $\max_{x \in C} \langle a, x \rangle$*
- (b) *c solves $\max_{x \in C} \langle a, x \rangle - \frac{1}{k} \|x - c\|^2$.*

(ii) *The problem in (b) reduces to a projection problem:*

$$\arg \max_{x \in C} \langle a, x \rangle - \frac{1}{k} \|x - c\|^2 = \arg \min_{x \in C} \|x - c - \frac{k}{2} a\|^2. \quad (4.1)$$

The following proposition characterizes projection on a unit simplex.

PROPOSITION 4.2. *Let $P : \mathbb{R}^n \rightarrow \Delta_n$ denote the projection on the $(n - 1)$ -dimensional unit simplex Δ_n with respect to the standard Euclidean distance.*

- (i) P is Lipschitz continuous.
(ii) For every $x \in \mathbb{R}^n$ the projection on Δ_n can be rewritten as follows

$$P(x) = ([x_1 + \lambda(x)]_+, \dots, [x_n + \lambda(x)]_+), \quad (4.2)$$

where $\lambda(x) \in \mathbb{R}$ is the unique solution to $\sum_{i=1}^n [x_i + \lambda(x)]_+ = 1$.

REMARK 4.3. Proposition 4.2(ii) immediately implies that for all $x \in \mathbb{R}^n$ and $i, j \in \{1, \dots, n\}$: if $x_i - x_j \geq 1$, then $P_j(x) = 0$.

4.3 The target projection dynamic

The target projection dynamic, the dynamic process governing the learning mechanism that we study in this paper, was mentioned briefly in the concluding section of Sandholm (2005, pp. 166-167). It was originally defined in the framework of congestion networks by Friesz et al. (1994).

DEFINITION 4.4. Consider a normal form game $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$. The **target projection dynamic (TPD)** is defined, for each $i \in N$, by the differential equation

$$\dot{\alpha}_i = P_{\Delta_i}[\alpha_i + U_i(\alpha)] - \alpha_i, \quad (4.3)$$

where P_{Δ_i} denotes the projection on Δ_i with respect to the usual Euclidean distance.

The basic idea is simple and standard for most dynamic processes in game theory. The payoffs associated with the different actions determine the direction in which their weights are changed by reinforcing the better actions and decreasing the weight of worse ones. Of course, simply running in the direction of the payoff vector $U_i(\alpha)$ might take you outside the strategy simplex, but Proposition 4.2(ii) assures that projection onto the strategy simplex does not affect the order of the coordinates.

The simplest dynamic embodying this idea is the best-response dynamic (Gilboa and Matsui, 1991): players place higher weight to their best responses according to

$$\dot{\alpha}_i = BR_i(\alpha) - \alpha_i, \quad (4.4)$$

where $BR_i(\alpha) = \arg \max_{\beta_i \in \Delta_i} \langle \beta_i, U_i(\alpha) \rangle$. In fact, Equation (4.4) is a differential inclusion, rather than a differential equation, since $BR_i(\alpha)$ need not be a singleton. In any case, as it becomes obvious from Equation (4.4), players do not switch to – but towards – their best responses. That is, the best-response dynamic involves inertia which precludes them from switching directly to their best response.

The next proposition indicates that the TPD is essentially a best-response dynamic, albeit under a bounded rationality assumption, involving the introduction of a certain status-quo bias. We follow the control cost approach, which since its introduction by Van Damme (1991) in the study of equilibrium refinements has proved to be a versatile way of providing microeconomic foundations for a variety of models of strategic behavior (Hofbauer and Sandholm, 2002; Mattsson and Weibull, 2002; Voorneveld, 2006). It does so by showing that such behavior is rational for decision makers who have to make some effort (incur costs) to implement their strategic choices. One intuitive way of modeling status-quo bias by player i could be as follows. Suppose that deviation from the current α_i is costly/requires effort in the sense that by switching to a strategy β_i , player i incurs a cost of $\frac{1}{2} \|\beta_i - \alpha_i\|^2$: staying at the current mixed strategy is costless, whereas large deviations, i.e., strategies further away from the current one in terms of Euclidean distance, incur larger costs. Taking such costs into account changes the optimization problem to

$$\max_{\beta_i \in \Delta_i} \langle \beta_i, U_i(\alpha) \rangle - \frac{1}{2} \|\beta_i - \alpha_i\|^2. \quad (4.5)$$

Let $B_i(\alpha) \in \Delta_i$ denote player i 's (unique due to strict concavity of the goal function) best response against α , i.e., the unique solution to problem (4.5). Subject to these assumptions, we can now formulate the TPD as a best response dynamic:

PROPOSITION 4.5. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game. The TPD is the best response dynamic for the control cost problem in (4.5), i.e., for each $i \in N$:*

$$\dot{\alpha}_i = P_{\Delta_i}[\alpha_i + U_i(\alpha_{-i})] - \alpha_i = B_i(\alpha) - \alpha_i. \quad (4.6)$$

PROOF. By definition,

$$P_{\Delta_i}[\alpha_i + U_i(\alpha)] = \arg \min_{\beta_i \in \Delta_i} \|\beta_i - \alpha_i - U_i(\alpha)\|^2, \quad (4.7)$$

so we need to establish that

$$\arg \min_{\beta_i \in \Delta_i} \|\beta_i - \alpha_i - U_i(\alpha)\|^2 = \arg \max_{\beta_i \in \Delta_i} \langle \beta_i, U_i(\alpha) \rangle - \frac{1}{2} \|\beta_i - \alpha_i\|^2,$$

which follows from Proposition 4.1(ii) for $C = \Delta_i$, $x = \beta_i$, $c = \alpha_i$, $a = U_i(\alpha)$ and $k = 2$. *QED.*

Proposition 4.5 shows that the TPD is subject to two types. Firstly, it preserves the inertia that all best-response dynamics exhibit, i.e., the players do not switch to any of their best responses, but they place higher weight to them, thus shifting their behavior towards them. Secondly, players dislike moving away from their current strategies. The fact that they are averse to changing their current behavior makes the TPD a conservative rule of adjustment, and learning becomes slow.

As pointed out to us by Bill Sandholm, the TPD actually belongs to a larger family of dynamics, characterized by the degree of aversion towards shifting away from the current strategy. Consider the parametric dynamic

$$\dot{\alpha}_i = P_{\Delta_i}[\alpha_i + \frac{k}{2}U_i(\alpha)] - \alpha_i, \quad (4.8)$$

where $k > 0$. Larger k induces (a) higher weight placed on the strategic component of the adjustment process, i.e., on moving towards the best response, and (b) lower weight to the cost incurred by changing the current behavior. This becomes clearer by looking at Proposition 4.1(ii): increasing k leads to lower control costs, since the weight placed on these costs by the perturbed payoff function is equal to $\frac{1}{k}$. The TPD arises when $k = 2$. When $k \rightarrow \infty$ the control cost becomes arbitrarily small and (4.8) converges to a (unique and Lipschitz continuous) solution trajectory of the best-response dynamic (see Theorem 4.6 in the following section). On the other hand, the lower k the more conservative the players become, since they start placing more weight on their aversion towards change. In the limit ($k \rightarrow 0$), the players never change their actual strategy and stick forever to it. In the present paper we focus on the TPD ($k = 2$), but it is worth mentioning that all our results hold for every $k > 0$, since (4.8) corresponds to the TPD for a different game. That is, the parametric family of TPD in (4.8), no matter how conservative rules of adjustment it imposes, has a series of nice properties as shown in the following sections.

The fact that the TPD is interpreted as a best-response dynamic with control costs does not imply that the dynamic is susceptible to the extensive literature on (perturbed) best

response dynamics (Gilboa and Matsui, 1991; Hofbauer and Sandholm, 2002; Fudenberg and Levine, 1998): Our status-quo bias models control costs arising due to deviations from the *current* state, while the usual approaches use control cost functions which:

- are independent of the current state: they define costs in terms of deviations from a fixed strategy, often uniform randomization (close your eyes and pick an action) as in Mattsson and Weibull (2002), and Voorneveld (2006),
- are often required to be steep near the boundary of the strategy space, as in Hofbauer and Sandholm (2002).

4.3.1 General properties

Theorem 4.6 states that the target projection dynamic satisfies a number of desirable properties of “nice” evolutionary dynamics. Indeed, Sandholm (2005) calls a dynamic *well-behaved* if it satisfies the first three properties of Theorem 4.6.

THEOREM 4.6. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game. The TPD satisfies the following properties:*

NASH STATIONARITY: *The stationary points of the TPD and the game’s Nash equilibria coincide.*

BASIC SOLVABILITY: *For every initial state, a solution to the TPD exists, is unique, Lipschitz continuous in the initial state, and remains inside Δ at all times.*

POSITIVE CORRELATION: *Growth rates are positively correlated with payoffs: for each $i \in N$, if $\dot{\alpha}_i \neq 0$, then $\langle \dot{\alpha}_i, U_i(\alpha) \rangle > 0$.*

INNOVATION: *If some player is not at a stationary state and has an unused best response, then a positive probability is assigned to it. Formally, for each $\alpha \in \Delta$ and $i \in N$, if $\dot{\alpha}_i \neq 0$, but there is an action $a_i^j \in A_i$ with $U_i^j(\alpha) = \max_{k \in \{1, \dots, J_i\}} U_i^k(\alpha)$ and $\alpha_i^j = 0$, then $\dot{\alpha}_i^j > 0$.*

PROOF. NASH STATIONARITY: Let $\alpha \in \Delta$. By Proposition 4.1, Proposition 4.5, and (4.3), the following chain of equivalences holds:

$$\alpha \text{ is a Nash equilibrium} \Leftrightarrow \forall i \in N : \alpha_i \in \arg \max_{\beta_i \in \Delta_i} \langle \beta_i, U_i(\alpha) \rangle$$

$$\begin{aligned}
&\Leftrightarrow \forall i \in N : \alpha_i \in \arg \max_{\beta_i \in \Delta_i} \langle \beta_i, U_i(\alpha) \rangle - \frac{1}{2} \|\beta_i - \alpha_i\|^2 \\
&\Leftrightarrow \forall i \in N : \alpha_i = P_{\Delta_i}[\alpha_i + U_i(\alpha)] \\
&\Leftrightarrow \forall i \in N : \dot{\alpha}_i = 0.
\end{aligned}$$

BASIC SOLVABILITY: The target projection dynamic (4.3) is Lipschitz continuous. Let $i \in N$. By assumption, the payoff U_i is Lipschitz continuous, say with expansion factor $C > 0$. By Proposition 4.2, the projection is Lipschitz continuous with expansion factor 1. Using the triangle inequality, it follows for each $\alpha, \beta \in \Delta$ that

$$\begin{aligned}
\|P[\alpha_i + U_i(\alpha)] - \alpha_i - P[\beta_i + U_i(\beta)] + \beta_i\| &\leq \|P[\alpha_i + U_i(\alpha)] - P[\beta_i + U_i(\beta)]\| + \|\alpha_i - \beta_i\| \\
&\leq \|\alpha_i + U_i(\alpha) - \beta_i - U_i(\beta)\| + \|\alpha_i - \beta_i\| \\
&\leq \|U_i(\alpha) - U_i(\beta)\| + 2\|\alpha_i - \beta_i\| \\
&\leq (C + 2)\|\alpha_i - \beta_i\|,
\end{aligned}$$

establishing Lipschitz continuity of the vector field in (4.3). Since P_{Δ_i} maps onto Δ_i , it follows that $\sum_{j=1}^{J_i} \dot{\alpha}_i^j = 0$. Moreover, if $\alpha_i^j = 0$, then $\dot{\alpha}_i^j \geq 0$. This makes Δ forward-invariant. Together, these properties imply (Hirsch and Smale, 1974, Ch. 8) that for every initial state, a solution exists, is unique, Lipschitz continuous in the initial state, and remains in Δ at all times.

POSITIVE CORRELATION: Let $\alpha \in \Delta$ and $i \in N$. Suppose $\dot{\alpha}_i = P_{\Delta_i}[\alpha_i + U_i(\alpha)] - \alpha_i \neq 0$. Let $\beta_i = P_{\Delta_i}[\alpha_i + U_i(\alpha)] \neq \alpha_i$. Then, using Proposition 4.5, one obtains:

$$\begin{aligned}
\langle \beta_i, U_i(\alpha) \rangle &> \langle \beta_i, U_i(\alpha) \rangle - \frac{1}{2} \|\alpha_i - \beta_i\|^2 \\
&\geq \langle \alpha_i, U_i(\alpha) \rangle - \frac{1}{2} \|\alpha_i - \alpha_i\|^2 \\
&= \langle \alpha_i, U_i(\alpha) \rangle,
\end{aligned}$$

So $\langle \dot{\alpha}_i, U_i(\alpha) \rangle = \langle \beta_i - \alpha_i, U_i(\alpha) \rangle > 0$.

INNOVATION: Assume that the premises of the innovation property hold, but that $\dot{\alpha}_i^j \leq 0$. We derive a contradiction. By Proposition 4.2 there is a $\lambda \in \mathbb{R}$ such that

$$P_{\Delta_i}[\alpha_i + U_i(\alpha)] = ([\alpha_i^1 + U_i^1(\alpha) + \lambda]_+, \dots, [\alpha_i^{J_i} + U_i^{J_i}(\alpha) + \lambda]_+).$$

By assumption, action j is unused ($\alpha_i^j = 0$) and $\dot{\alpha}_i^j \leq 0$, so

$$0 \geq \dot{\alpha}_i^j = [\alpha_i^j + U_i^j(\alpha) + \lambda]_+ - \alpha_i^j = [U_i^j(\alpha) + \lambda]_+ \geq 0,$$

i.e., $\dot{\alpha}_i^j = 0$ and $U_i^j(\alpha) + \lambda \leq 0$. But action j is a best response: $U_i^j(\alpha) = \max_{k \in \{1, \dots, J_i\}} U_i^k(\alpha)$. Consequently, for every action $k \in \{1, \dots, J_i\}$:

$$\dot{\alpha}_i^k = [\alpha_i^k + U_i^k(\alpha) + \lambda]_+ - \alpha_i^k \leq [\alpha_i^k + U_i^j(\alpha) + \lambda]_+ - \alpha_i^k \leq [\alpha_i^k + 0]_+ - \alpha_i^k = 0.$$

Since $\sum_{k=1}^{J_i} \dot{\alpha}_i^k = 0$, this implies that $\dot{\alpha}_i^k = 0$ for all $k \in \{1, \dots, J_i\}$, in contradiction with the assumption that $\dot{\alpha}_i \neq 0$. *QED.*

4.3.2 Strict domination: mind the gap

Berger and Hofbauer (2006) show that under the Brown-von Neumann-Nash (BNN) dynamic, introduced in Brown and von Neumann (1950), there are games where a strictly dominated strategy survives. Hofbauer and Sandholm (2007) generalize this example: for each evolutionary dynamic satisfying the properties in Theorem 4.6 –actually, they restrict attention to single-population games– it is possible to construct a game with a strictly dominated strategy that survives along solutions of most initial states.

As their result applies to our TPD, it is of interest to investigate whether there are additional conditions under which such “bad” actions *are* wiped out. The next result shows that this is the case if one action strictly dominates another and the “gap” between them is sufficiently large.

PROPOSITION 4.7. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game and let $i \in N$. Suppose there are actions $k, \ell \in \{1, \dots, J_i\}$ such that $U_i^k - U_i^\ell \geq 2$, i.e., action k strictly dominates action ℓ , and the gap between the payoffs is at least two. Then the probability α_i^ℓ converges to zero in the TPD.*

PROOF. We show that the differential equation for the probability α_i^ℓ of action ℓ is given by $\dot{\alpha}_i^\ell = -\alpha_i^\ell$, because then $\alpha_i^\ell(t) = \alpha_i^\ell(0)e^{-t} \rightarrow 0$ as $t \rightarrow \infty$. Let $\alpha \in \Delta$. By (4.3), it suffices to show that the ℓ -th coordinate of the projection $P_{\Delta_i}[\alpha_i + U_i(\alpha)]$ is zero. By Proposition 4.2(ii), there is a $\lambda \in \mathbb{R}$ such that its ℓ -th and k -th coordinate can be written as $[\alpha_i^\ell + U_i^\ell(\alpha) + \lambda]_+$ and $[\alpha_i^k + U_i^k(\alpha) + \lambda]_+$. Suppose, contrary to what we want to prove, that $[\alpha_i^\ell + U_i^\ell(\alpha) + \lambda]_+ > 0$. Then

$$\begin{aligned}
[\alpha_i^k + U_i^k(\alpha) + \lambda]_+ - [\alpha_i^\ell + U_i^\ell(\alpha) + \lambda]_+ &\geq \alpha_i^k - \alpha_i^\ell + U_i^k(\alpha) - U_i^\ell(\alpha) \\
&\geq \alpha_i^k - \alpha_i^\ell + 2 \\
&\geq 1,
\end{aligned} \tag{4.9}$$

since the difference between probabilities is bounded in absolute value by one. However, since $[\alpha_i^\ell + U_i^\ell(\alpha) + \lambda]_+ > 0$, the left-hand side of (4.9) is smaller than one, a contradiction. *QED.*

Also if a player has only two actions to choose from, and one of them is strictly dominated, then it is eventually eliminated:

PROPOSITION 4.8. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game and let $i \in N$. If $A_i = \{a_i^1, a_i^2\}$, and a_i^1 strictly dominates a_i^2 , the probability assigned to a_i^2 converges to zero in the TPD.*

PROOF. Let $\alpha \in \Delta$. By Proposition 4.2, there is a $\lambda(\alpha) \in \mathbb{R}$ such that the target projection dynamic for each of the two actions $j = 1, 2$ of population i can be rewritten as

$$\dot{\alpha}_i^j = [\alpha_i^j + U_i^j(\alpha) + \lambda(\alpha)]_+ - \alpha_i^j \tag{4.10}$$

Then $[\alpha_i^1 + U_i^1(\alpha) + \lambda(\alpha)]_+ > 0$. Suppose, to the contrary, that

$$\alpha_i^1 + U_i^1(\alpha) + \lambda(\alpha) \leq 0.$$

Since we project the two-dimensional vector onto the simplex, this implies

$$[\alpha_i^2 + U_i^2(\alpha) + \lambda(\alpha)]_+ = \alpha_i^2 + U_i^2(\alpha) + \lambda(\alpha) = 1.$$

Combining these two expressions gives

$$\alpha_i^2 - \alpha_i^1 \geq 1 + U_i^1(\alpha) - U_i^2(\alpha) > 1,$$

a contradiction, since the left-hand side is at most one. By continuity of the payoffs on the compact set Δ and strict domination, there is an $\varepsilon > 0$ such that $U_i^1(\alpha) - U_i^2(\alpha) > \varepsilon$ for each $\alpha \in \Delta$.

Distinguish two cases. First, if $[\alpha_i^2 + U_i^2(\alpha) + \lambda(\alpha)]_+ = 0$, then $\dot{\alpha}_i^2 = -\alpha_i^2$, so the probability α_i^2 decreases at an exponential rate. Second, if $[\alpha_i^2 + U_i^2(\alpha) + \lambda(\alpha)]_+ > 0$, combine this with the facts that $[\alpha_i^1 + U_i^1(\alpha) + \lambda(\alpha)]_+ > 0$ and that these two numbers add up to one, to deduce that $\lambda(\alpha) = -\frac{1}{2}(U_i^1(\alpha) + U_i^2(\alpha))$. So $\dot{\alpha}_i^2 = \frac{1}{2}(U_i^2(\alpha) - U_i^1(\alpha)) < -\frac{1}{2}\varepsilon$, i.e., the probability α_i^2 decreases at a rate bounded away from zero. Hence, along any solution trajectory, the probability α_i^2 of the dominated action converges to zero. *QED.*

4.3.3 The projection dynamic and the target projection dynamic

The projection dynamic was first developed by Nagurney and Zang (1997) as part of the transportation literature, and was later introduced to game theory by Sandholm (2006), Lahkar and Sandholm (2008), and Sandholm et al. (2008). Let $T_i = \{\beta_i \in \mathbb{R}^{J_i} : \sum_{j=1}^{J_i} \beta_i^j = 0\}$ be the tangent space of Δ_i . Every $z_i \in T_i$ describes a motion between two points in Δ_i . The tangent cone of Δ_i at some strategy $\alpha_i \in \Delta_i$ is the set of feasible motions from α_i towards some other strategy in Δ_i , i.e., $T_i(\alpha) = \{\beta_i \in T_i : \alpha_i^j = 0 \Rightarrow \beta_i^j \geq 0\}$. Then, the projection dynamic system is defined as follows

$$\dot{\alpha}_i = P_{T_i(\alpha)}[U_i(\alpha)], \quad (4.11)$$

The following proposition shows that the projection dynamic and the target projection dynamic coincide at α if α is completely mixed and the target projection is orthogonal to the motion it causes.

PROPOSITION 4.9. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game. If $\alpha \in \text{int}(\Delta)$ and $\langle \dot{\alpha}_i, U_i(\alpha) - \dot{\alpha}_i \rangle = 0$, with $\dot{\alpha}$ as in (4.3), the TPD coincides with the projection dynamic.*

PROOF. By definition the TPD solves, for all $i \in \mathbb{N}$, the following optimization problem:

$$\min_{\beta_i \in \Delta_i} \|\beta_i - \alpha_i - U_i(\alpha)\|^2, \quad \text{s.t.} \quad \sum_{j=1}^{J_i} \beta_i^j = 1 \text{ and } \beta_i^j \geq 0 \text{ for all } j = 1, \dots, J_i.$$

It follows from the Karush-Kuhn-Tucker conditions that there are $\nu \in \mathbb{R}$ and $\mu_j \geq 0$ such that for all $j = 1, \dots, J_i$

$$\begin{aligned} \beta_i^j - \alpha_i^j - U_i^j(\alpha) + \nu - \mu_j &= 0, \\ \mu_j \beta_i^j &= 0. \end{aligned} \quad (4.12)$$

Summing (4.12) for all j , solving with respect to ν and substituting back yields

$$\beta_i^j - \alpha_i^j - U_i^j(\alpha) + \frac{1}{J_i} \sum_{j=1}^{J_i} U_i^j(\alpha) - \mu_j + \frac{1}{J_i} \sum_{j=1}^{J_i} \mu_j = 0. \quad (4.13)$$

We multiply by β_i^j , sum for all j and use the complementary slackness condition ($\mu_j \beta_i^j = 0$):

$$\sum_{j=1}^{J_i} \left(\beta_i^{j^2} - \alpha_i^j \beta_i^j - \beta_i^j U_i^j(\alpha) \right) + \frac{1}{J_i} \sum_{j=1}^{J_i} U_i^j(\alpha) + \frac{1}{J_i} \sum_{j=1}^{J_i} \mu_j = 0. \quad (4.14)$$

We multiply now (4.13) by α_i^j and sum for all j :

$$\sum_{j=1}^{J_i} \left(\alpha_i^j \beta_i^j - \alpha_i^{j^2} - \alpha_i^j U_i^j(\alpha) \right) + \frac{1}{J_i} \sum_{j=1}^{J_i} U_i^j(\alpha) - \sum_{j=1}^{J_i} \mu_j \alpha_i^j + \frac{1}{J_i} \sum_{j=1}^{J_i} \mu_j = 0. \quad (4.15)$$

Now subtract (4.14) from (4.15) to find

$$\langle \beta_i - \alpha_i, \beta_i - \alpha_i - U_i(\alpha) \rangle = \sum_{j=1}^{J_i} \mu_j \alpha_i^j.$$

As $\beta_i = P_{\Delta_i}[\alpha_i + U_i(\alpha)]$, the orthogonality assumption implies that $\langle \beta_i - \alpha_i, \beta_i - \alpha_i - U_i(\alpha) \rangle = 0$. As $\mu_j \alpha_i^j \geq 0$ for all $j = 1, \dots, J_i$ and $\alpha \in \text{int}(\Delta)$, it follows that $\mu_j = 0$ for all pure strategies $j = 1, \dots, J_i$. Then, it follows from Proposition 4.2(ii) that the projection can be rewritten as

$$P_{\Delta_i}[\alpha_i + U_i(\alpha)] = ([\alpha_i^1 + U_i^1(\alpha) + \lambda]_+, \dots, [\alpha_i^{J_i} + U_i^{J_i}(\alpha) + \lambda]_+),$$

with $\alpha_i^j + U_i^j(\alpha) + \lambda \geq 0$ for all j , since $\mu_j = 0$. Hence, $\lambda = -\frac{1}{J_i} \sum_{j=1}^{J_i} U_i^j(\alpha)$, which implies that for every $i \in N$ and every $j \in J_i$ the TPD becomes

$$\dot{\alpha}_i^j = U_i^j(\alpha) - \frac{1}{J_i} \sum_{k=1}^{J_i} U_i^k(\alpha).$$

The previous formula is the projection dynamic for all completely mixed strategies (Sandholm et al., 2008), which proves the proposition. *QED.*

COROLLARY 4.10. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game and let α be a completely mixed Nash equilibrium. Then, there is a neighborhood \mathcal{O} of α such that the projection dynamic and the target projection dynamic coincide for all $\beta \in \mathcal{O}$.*

PROOF. It follows directly from Proposition 4.9 and continuity in \mathcal{O} . *QED.*

4.4 Special classes of games

In this section we study the properties of the TPD in some special classes of games.

4.4.1 Stable games

Sandholm et al. (2008) prove a number of stability results for potential and stable games under the projection dynamic. Stable games (Hofbauer and Sandholm, 2008) are a family of normal form games characterized by the following condition:

$$\langle \alpha_i - \beta_i, U_i(\alpha) - U_i(\beta) \rangle \leq 0, \quad (4.16)$$

for every $\alpha_i, \beta_i \in \Delta_i$ and for all $i \in N$. The game is null (strictly) stable if (4.16) holds with equality (strict inequality).

PROPOSITION 4.11. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game and let α be a completely mixed Nash equilibrium.*

- (i) *If the game is stable then α is Lyapunov stable under the TPD.*
- (ii) *If the game is strictly stable then α is asymptotically stable under the TPD.*
- (iii) *If the game null stable then there is a neighborhood of α where the squared Euclidean distance to α , i.e., the function $\beta \mapsto \|\beta - \alpha\|^2$, defines a constant of motion under the TPD.*

PROOF. By Corollary 4.10, there is a neighborhood \mathcal{O} of α where the projection dynamic agrees with the target projection dynamic for all players. Let $\varepsilon > 0$ be such that

$$\mathcal{O}_\varepsilon := \{\beta \in \Delta : \|\beta - \alpha\|^2 \leq \varepsilon\} \subseteq \mathcal{O}.$$

Sandholm et al. (2008) show that in the three cases of our proposition the function $L : \Delta \rightarrow \mathbb{R}$ with $L(\beta) := \|\beta - \alpha\|^2$ is (i) a Lyapunov function, (ii) a strict Lyapunov function, and (iii) defines a constant motion around α under the projection dynamic. Since $\mathcal{O}_\varepsilon \subseteq \mathcal{O}$, every trajectory starting in \mathcal{O}_ε will remain in it forever under the projection dynamic, and therefore under the target projection dynamic, completing the proof. *QED.*

4.4.2 Zero-sum games

A very interesting and widely explored class of games is the zero-sum games. We say that a normal form game is zero-sum if $\sum_{i \in N} u_i(a) = 0$ for every $a \in A$. Hofbauer and Sandholm (2008) establish that zero-sum games are null stable. By Proposition 4.11(iii), every trajectory of the TPD that gets sufficiently close to a completely mixed equilibrium in two-players zero-sum games forms a closed cyclical orbit around it.

COROLLARY 4.12. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a two-player normal form zero-sum game. If α is a completely mixed Nash equilibrium, there is a neighborhood \mathcal{O} of α on which the TPD forms a constant of motion around α .*

That is, if $\beta_0 \in \mathcal{O}$ and β belongs to the trajectory of the (unique) solution of (4.3) with initial value β_0 , then $\|\beta - \alpha\| = \|\beta_0 - \alpha\|$. Typical examples include the matching pennies and the rock-paper-scissors games.

4.4.3 Games with strict Nash equilibria

Recall that in finite strategic games a Nash equilibrium is strict if each player chooses the unique best reply, i.e., $\alpha \in \Delta$ is a strict Nash equilibrium if $\langle \alpha_i, U_i(\alpha) \rangle > \langle \beta_i, U_i(\alpha) \rangle$ for all $\beta \in \Delta$ and all $i \in N$. Consequently, strict Nash equilibria are equilibria in pure strategies. This follows from the fact that a mixed strategy α_i is a best response to α if and only if all actions assigned positive probability by α_i are best responses to α , implying that i is indifferent between these actions, and therefore the necessary and sufficient condition for the strict equilibrium is violated.

PROPOSITION 4.13. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a normal form game. If α is a strict Nash equilibrium, it is asymptotically stable under the TPD.*

PROOF. Let α be a strict Nash equilibrium. Since α must be in pure strategies, without loss of generality, each $i \in N$ plays his first action: $\alpha_i = a_i^1$. By definition, for each $i \in N$ and $j \in \{2, \dots, J_i\}$: $U_i^1(\alpha) > U_i^j(\alpha)$, so that $(\alpha_i^1 + U_i^1(\alpha)) - (\alpha_i^j + U_i^j(\alpha)) = 1 + U_i^1(\alpha) - U_i^j(\alpha) > 1$. By continuity, there is a neighborhood \mathcal{O} of α such that for all $\beta \in \mathcal{O}$, $i \in N$, and $j \in \{2, \dots, J_i\}$:

$$(\beta_i^1 + U_i^1(\beta)) - (\beta_i^j + U_i^j(\beta)) \geq 1.$$

For all $\beta \in \mathcal{O}$ and $i \in N$, Remark 4.3 implies that $P_{\Delta_i}(\beta_i + U_i(\beta)) = \alpha_i$; so $\dot{\beta}_i = \alpha_i - \beta_i$. Hence, the function $L : \mathcal{O} \rightarrow \mathbb{R}$ with $L(\beta) := \sum_{i \in N} \|\beta_i - \alpha_i\|^2$ is a Lyapunov function: It is non-negative, zero only at α , and if $\beta \in \mathcal{O} \setminus \{\alpha\}$:

$$\dot{L} = 2 \sum_{i \in N} \langle \beta_i - \alpha_i, \dot{\beta}_i \rangle = 2 \sum_{i \in N} \langle \beta_i - \alpha_i, \alpha_i - \beta_i \rangle = -2 \sum_{i \in N} \|\beta_i - \alpha_i\|^2 < 0.$$

Given the existence of the Lyapunov function L , the equilibrium α is asymptotically stable (Hirsch and Smale, 1974, Ch. 9). *QED.*

4.4.4 Games with evolutionarily stable strategies

We focus on symmetric two-player normal form games. A two player game is called symmetric if $A_1 = A_2 = A$ and $u_1(a_1, a_2) = u_2(a_1, a_2) = u(a_1, a_2)$ for all $a_1, a_2 \in A$. We say that (α, α) is an ***evolutionarily stable strategy (ESS)*** of a symmetric two-player normal form game (Maynard Smith, 1982; Weibull, 1995; Fudenberg and Levine, 1998) if for all $\beta \neq \alpha$:

- (a) $\langle \alpha, U(\alpha) \rangle > \langle \beta, U(\alpha) \rangle$, or
- (b) $\langle \alpha, U(\alpha) \rangle = \langle \beta, U(\alpha) \rangle$ and $\langle \alpha, U(\beta) \rangle > \langle \beta, U(\beta) \rangle$.

That is, a strategy is evolutionarily stable if it is robust to behavioral mutations: small mutations receive a strictly lower post-entry payoff than the incumbent strategy. ESS is a refinement of Nash equilibrium, so all ESS are rest points of every dynamic that satisfies Nash stationarity. Since the definition of evolutionary stability is conceptually based on the idea that the mutants eventually assimilate to the original population, one would expect ESS to attract trajectories that get sufficiently close to them under dynamic processes of myopic adjustment. Taylor and Jonker (1978), and Hofbauer et al. (1979) show that every ESS is asymptotically stable under the replicator dynamic. However, a similar result cannot be established for the TPD. Instead we provide some partial stability results.

Games with a completely mixed evolutionary stable strategy are strictly stable (Hofbauer and Sandholm, 2008). Proposition 4.11(ii) thus implies:

COROLLARY 4.14. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a symmetric two-player normal form game and let α be a completely mixed ESS. Then, it is asymptotically stable in the TPD.*

If in addition we restrict players to choose between only two actions we can extend the previous result to all ESS:

PROPOSITION 4.15. *Let $(N, (A_i)_{i \in N}, (U_i)_{i \in N})$ be a symmetric 2×2 normal form game and let α be an ESS. Then, it is asymptotically stable in the TPD.*

PROOF. For convenience, denote the strategy space by $\Delta = \{\beta \in \mathbb{R}_+^2 : \beta_1 + \beta_2 = 1\}$. Let $A = \{a_1, a_2\}$ and let $\alpha \in \Delta$ be an ESS. If $\alpha = (\alpha_1, \alpha_2)$ is completely mixed, apply Corollary 4.14. If not, assume without loss of generality that $\alpha = e_1$, where $e_1 = (1, 0)$, i.e., the pure strategy a_1 is an ESS.

If $U_1(e_1) > U_2(e_1)$ then it follows from convexity of the payoffs that

$$\langle e_1, U(e_1) \rangle = U_1(e_1) > \beta_1 U_1(e_1) + \beta_2 U_2(e_1) = \langle \beta, U(e_1) \rangle,$$

for all $\beta = (\beta_1, \beta_2) \in \Delta$. Hence, from Proposition 4.13 it follows that e_1 is asymptotically stable in the TPD.

Suppose now that $U_1(e_1) = U_2(e_1)$. Since e_1 is an ESS, it is also a Nash equilibrium, and therefore due to Theorem 4.6 it is a rest point. This implies that $P_\Delta[e_1^i + U_i(e_1)] = e_1^i$, which yields $P_\Delta[1 + U_1(e_1)] = 1 > 0$ for $i = 1$. It follows from continuity that there is a neighborhood \mathcal{O} of e_1 such that $P_\Delta[\beta_1 + U_1(\beta)] > 0$ for all $\beta \in \mathcal{O}$. From Proposition 4.2(ii) it follows that there is $\lambda(\beta) \in \mathbb{R}$ such that $P_\Delta[\beta_1 + U_1(\beta)] = [\beta_1 + U_1(\beta) + \lambda(\beta)]_+ > 0$, implying that $P_\Delta[\beta_1 + U_1(\beta)] = \beta_1 + U_1(\beta) + \lambda(\beta)$, for all $\beta \in \mathcal{O}$. Hence,

$$\dot{\beta}_1 = U_1(\beta) + \lambda(\beta), \tag{4.17}$$

for all $\beta \in \mathcal{O}$. Now, we consider two cases:

CASE 1: Let $[\beta_2 + U_2(\beta) + \lambda(\beta)]_+ > 0$, which implies again due to Proposition 4.2(ii) that $\lambda(\beta) = -\frac{1}{2}(U_1(\beta) + U_2(\beta))$. Hence, it follows from Equation (4.17) that

$$\begin{aligned} \dot{\beta}_1 &= U_1(\beta) - \frac{1}{2}(U_1(\beta) + U_2(\beta)) \\ &= \frac{1}{2}(U_1(\beta) - U_2(\beta)) \\ &= \frac{1}{2}\beta_1 \left(\langle e_1, U_1(e_1) \rangle - \langle e_2, U_2(e_1) \rangle \right) + \frac{1}{2}\beta_2 \left(\langle e_1, U_1(e_2) \rangle - \langle e_2, U_2(e_2) \rangle \right) \\ &= \frac{1}{2}\beta_2 \left(\langle e_1, U_1(e_2) \rangle - \langle e_2, U_2(e_2) \rangle \right) > 0, \end{aligned} \tag{4.18}$$

where $e_2 = (0, 1)$ denotes the pure strategy a_2 .

CASE 2: Let $[\beta_2 + U_2(\beta) + \lambda(\beta)]_+ = 0$, which implies again due to Proposition 4.2(ii) that $\lambda(\beta) = 1 - \beta_1 - U_1(\beta)$. Substituting into Equation (4.17) yields

$$\dot{\beta}_1 = 1 - \beta_1 > 0. \quad (4.19)$$

Consider now $L(\beta) = 1 - \beta_1$ for $\beta \in \mathcal{O}$, which is a Lyapunov function: it is non-negative, equal to zero only at e_1 , and for all $\beta \in \mathcal{O} \setminus \{e_1\}$

$$\dot{L} = -\dot{\beta}_1 < 0,$$

which follows from (4.18) and (4.19) and completes the proof. *QED.*

Appendix

PROOF OF PROPOSITION 4.1. (i) [(a) \Rightarrow (b)] Assume (a) holds. Since $\|c - c\| = 0$, it follows, for each $x \in C$, that

$$\langle a, c \rangle - \frac{1}{k}\|c - c\|^2 = \langle a, c \rangle \geq \langle a, x \rangle \geq \langle a, x \rangle - \frac{1}{k}\|x - c\|^2,$$

so (b) holds.

[(b) \Rightarrow (a)] Assume (b) holds. Let $x \in C$ and $\lambda \in (0, 1)$. By convexity, $\lambda x + (1 - \lambda)c \in C$. Since $\|c - c\| = 0$, it follows that

$$\begin{aligned} \langle a, c \rangle &\geq \langle a, \lambda x + (1 - \lambda)c \rangle - \frac{1}{k}\|(\lambda x + (1 - \lambda)c) - c\|^2 \\ &= \lambda \langle a, x \rangle + (1 - \lambda) \langle a, c \rangle - \frac{\lambda^2}{k}\|x - c\|^2. \end{aligned}$$

Rearrange terms and divide by $\lambda > 0$ to obtain that $\langle a, c \rangle \geq \langle a, x \rangle - \frac{\lambda}{k}\|x - c\|^2$. Since $\lambda \in (0, 1)$ is arbitrary, let λ approach zero to establish (a).

(ii) Maximizing the function $x \mapsto \langle a, x \rangle - \frac{1}{k}\|x - c\|^2$ is equivalent with minimizing $x \mapsto \frac{1}{k}\|x - c\|^2 - \langle a, x \rangle$. It therefore suffices to show that the latter function is a positive affine transformation of $x \mapsto \|x - c - \frac{k}{2}a\|^2$. Using the linearity and symmetry properties of the inner product, we find

$$\begin{aligned}
\|x - c - \frac{k}{2}a\|^2 &= \langle x - c - \frac{k}{2}a, x - c - \frac{k}{2}a \rangle \\
&= \langle x - c, x - c \rangle - 2\langle \frac{k}{2}a, x \rangle + 2\langle c, \frac{k}{2}a \rangle + \langle \frac{k}{2}a, \frac{k}{2}a \rangle \\
&= \|x - c\|^2 - k\langle a, x \rangle + k\langle c, a \rangle + \frac{k^2}{4}\|a\|^2,
\end{aligned}$$

which completes the proof, since the final two terms are independent of x , and $\|x - c\|^2 - k\langle a, x \rangle$ is a simple rescaling of $\frac{1}{k}\|x - c\|^2 - \langle a, x \rangle$. *QED.*

PROOF OF PROPOSITION 4.2. (i) Recall from the Projection Theorem (see, for instance, Luenberger, 1969, p. 69) that for every $z \in \mathbb{R}^n$, $P(z)$ is characterized by $\langle z - P(z), w - P(z) \rangle \leq 0$ for all $w \in \Delta_n$. In particular, for all $x, y \in \mathbb{R}^n$:

$$\langle x - P(x), P(y) - P(x) \rangle \leq 0 \text{ and } \langle y - P(y), P(x) - P(y) \rangle \leq 0.$$

Write $\langle y - P(y), P(x) - P(y) \rangle = \langle P(y) - y, P(y) - P(x) \rangle$, add the two inequalities, and use Cauchy-Schwarz to establish

$$\begin{aligned}
0 &\geq \langle x - P(x) + P(y) - y, P(y) - P(x) \rangle \\
&= \|P(y) - P(x)\|^2 - \langle y - x, P(y) - P(x) \rangle \\
&\geq \|P(y) - P(x)\|^2 - \|y - x\| \|P(y) - P(x)\|.
\end{aligned}$$

Conclude that $\|P(y) - P(x)\| \leq \|y - x\|$, i.e., P is Lipschitz continuous with expansion factor 1.

(ii) Let $x \in \mathbb{R}^n$. The function $T : \mathbb{R} \rightarrow \mathbb{R}$ defined for each $\lambda \in \mathbb{R}$ by $T(\lambda) = \sum_{i=1}^n [x_i + \lambda]_+$ is the composition of continuous functions and therefore continuous itself. Let $m = \max\{x_1, \dots, x_n\}$. Then $T(\lambda) = 0$ for all $\lambda \in (-\infty, -m]$ and T is strictly increasing on $[-m, \infty)$, with $T(\lambda) \rightarrow \infty$ as $\lambda \rightarrow \infty$. By the Intermediate Value Theorem, there is a unique $\lambda(x) \in [-m, \infty)$ such that $T(\lambda(x)) = 1$.

By definition, $P(x)$ is the unique solution to $\min_{y \in \Delta_n} \frac{1}{2} \sum_{i=1}^n (y_i - x_i)^2$. This is a convex quadratic optimization problem with linear constraints, so the Karush-Kuhn-Tucker conditions are necessary and sufficient to characterize the minimum location: $y^* \in \Delta_n$ solves the problem if and only if there exist Lagrange multipliers $\mu_i \geq 0$ associated with the inequality constraints $y_i \geq 0$ and $\nu \in \mathbb{R}$ associated with the equality constraint $\sum_{i=1}^n y_i = 1$ such that for each $i = 1, \dots, n$:

$$y_i^* - x_i - \mu_i + \nu = 0, \quad (4.20)$$

$$\mu_i y_i^* = 0. \quad (4.21)$$

Condition (4.20) is the first order condition obtained from differentiating the Lagrangian

$$(y, \mu_1, \dots, \mu_n, \nu) \mapsto \frac{1}{2} \sum_{i=1}^n (y_i - x_i)^2 - \sum_{i=1}^n \mu_i y_i + \nu \left(\sum_{i=1}^n y_i - 1 \right)$$

with respect to y_i and condition (4.21) is the complementary slackness condition. It is now easy to see that $y^* := ([x_1 + \lambda(x)]_+, \dots, [x_n + \lambda(x)]_+)$ solves the minimization problem: set $\mu_i = 0$ if $[x_i + \lambda(x)]_+ > 0$, $\mu_i = -x_i - \lambda(x) \geq 0$ if $[x_i + \lambda(x)]_+ \leq 0$, and $\nu = -\lambda(x)$. Substitution in (4.20) and (4.21) shows that these necessary and sufficient conditions are satisfied. *QED.*

References

- AUMANN, R.J. (1976). Agreeing to disagree, *Annals of Statistics* 4, 1236–1239.
- (1997). Rationality and bounded rationality, *Games and Economic Behavior* 21, 2–14.
- (1999). Interactive epistemology I: knowledge, *International Journal of Game Theory* 28, 263–300.
- BACHARACH, M. (1985). Some extensions of a claim of Aumann in an axiomatic model of knowledge, *Journal of Economic Theory* 37, 167–190.
- BATTIGALLI, P., BONANNO, G. (1999). Recent results on belief, knowledge and the epistemic foundations of game theory, *Research in Economics* 53, 149–225.
- BERGER, U., HOFBAUER, J. (2006). Irrational behavior in the Brown-von Neumann-Nash dynamics, *Games and Economic Behavior* 56, 1–6.
- BRANDENBURGER, A., DEKEL, E. (1987). Common knowledge with probability 1, *Journal of Mathematical Economics* 16, 237–245.
- BRANDENBURGER, A., DEKEL, E., GEANAKOPOLOS, J. (1992). Correlated equilibrium with generalized information structures, *Games and Economic Behavior* 4, 182–201.
- BROWN, G., VON NEUMANN, J. (1950). Solutions of games by differential equations, *Annals of Mathematical Studies* 24, 73–79.
- BÖRGERS, T., SARIN, R. (1997). Learning through reinforcement and replicator dynamics, *Journal of Economic Theory* 77, 1–14.
- CHELLAS, B.F. (1980). *Modal logic: an introduction*, Cambridge University Press, Cambridge.

- DEKEL, E., LIPMAN, B., RUSTICHINI, A. (1998). Standard state-space models preclude unawareness, *Econometrica* 66, 159–173.
- CAVE, J.A.K. (1983). Learning to agree, *Economics Letters* 12, 147–152.
- FRIESZ, T., BERNSTEIN, D., MEHTA, N., TOBIN, R., GANJALIZADEH, S. (1994). Day-to-day dynamic network disequilibria and idealized traveler information systems, *Operations Research* 46, 1120–1136.
- FUDENBERG, D., LEVINE, D. (1998). *The Theory of Learning in Games*, MIT Press, Cambridge, Massachusetts.
- GEANAKOPOLOS, J. (1995). Common knowledge. *Handbook of Game Theory with Economic Applications*, edited by R.J. Aumann and S. Hart, Vol. II, Ch. 40, Elsevier, North-Holland.
- GEANAKOPOLOS, J., POLEMARCHAKIS, H., (1982). We can't disagree forever, *Journal of Economic Theory* 28, 192–200.
- GILBOA, I., MATSUI, A. (1991). Social stability and equilibrium, *Econometrica* 59, 859–867.
- HALPERN, J.Y. (2001). Alternative semantics for unawareness, *Games and Economic Behavior* 37, 321–339.
- HANSON, R. (1998). Consensus by identifying extremists, *Theory and Decision* 44, 293–301.
- (2002). Disagreement is unpredictable, *Economics Letters* 77, 365–369.
- HART, S., MAS-COLELL, A. (2003). Uncoupled dynamics do not lead to Nash equilibrium, *American Economic Review* 93, 1830–1836.
- HART, S., HEIFETZ, A., SAMET, D. (1996). “Knowing whether”, “knowing that”, and the cardinality of state spaces, *Journal of Economic Theory* 70, 249–256.
- HEIFETZ, A. (1996). Comment on consensus without common knowledge, *Journal of Economic Theory* 70, 273–277.
- HINTIKKA, J. (1962). *Knowledge and belief*, Cornell University Press, Ithaca, NY.
- HIRSCH, M., SMALE, S. (1974). *Differential equations, dynamical systems, and linear algebra*, New York: Academic Press.
- HOFBAUER, J., SCHUSTER, P., SIGMUND, K. (1979). A note on evolutionarily stable strategies and game dynamics, *Journal of Theoretical Biology* 81, 609–612.

- HOFBAUER, J., SANDHOLM, W. (2002). On the global convergence of stochastic fictitious play, *Econometrica* 70, 2265-2294.
- (2007). Survival of dominated strategies under evolutionary dynamics, Working paper, University of Wisconsin.
- (2008). Stable population games and integrability for evolutionary dynamics, Working paper, University of Wisconsin.
- HOPKINS, E. (2002). Two competing models of how people learn in games, *Econometrica* 70, 2141–2166.
- HOUY, N., MENAGER, L. (2007). Communication, consensus, and order: Who wants to speak first? *Journal of Economic Theory*, forthcoming.
- KRIPKE, S. (1963). Semantic analysis of modal logic, *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik* 9, 67–96.
- MCKELVEY, R.D., PAGE, T. (1986). Common knowledge, consensus, and aggregate information, *Econometrica* 54, 109–128.
- KOESSLER, F. (2001). Common knowledge and consensus with noisy communication, *Mathematical Social Sciences* 42, 139–159.
- KRASUCKI, P. (1996). Protocols forcing consensus, *Journal of Economic Theory* 70, 266–272.
- LAHKAR, R., SANDHOLM, W. (2008). The projection dynamic and the geometry of population games, *Games and Economic Behavior*, forthcoming.
- LIPMAN, B. (1995). Information processing and bounded rationality: A survey, *Canadian Journal of Economics* 28, 42–67.
- LUENBERGER, D. (1969). *Optimization by vector space methods*, New York: John Wiley & Sons.
- MAYNARD SMITH, J. (1982). *Evolution and the theory of games*, Cambridge University Press.
- MATTSSON, L.-G., WEIBULL, J. (2002). Probabilistic choice and procedurally bounded rationality, *Games and Economic Behavior* 41, 61–78.
- MILGROM, P. (1981). An axiomatic characterization of common knowledge, *Econometrica* 49, 219–222.

- MILGROM, P., STOKEY, N. (1980). Information, trade, and common knowledge, *Journal of Economic Theory* 26, 17–27.
- MODICA, S., RUSTICHINI, A. (1994). Awareness and partitional information structures, *Theory and Decision* 37, 107–124.
- (1999). Unawareness and partitional information structures, *Games and Economic Behavior* 27, 265–298.
- NAGURNEY, A., ZHANG, D. (1997). Projected dynamical systems in the formulation, stability analysis, and computation of fixed demand traffic network equilibria, *Transportation Science* 31, 147–158.
- NIELSEN, L.T. (1984). Common knowledge, communication, and convergence of beliefs, *Mathematical Social Sciences* 8, 1–14.
- (1995). Common knowledge of a multivariate aggregate statistic, *International Economic Review* 36, 207–216.
- NIELSEN, L.T., BRANDENBURGER, A., GEANAKOPOLOS, J., MCKELVEY, R., PAGE, T. (1990). Common knowledge of an aggregate of expectations, *Econometrica* 58, 1235–1239.
- PARIKH, R., KRASUCKI, P. (1990). Communication, consensus, and knowledge, *Journal of Economic Theory* 52, 178–189.
- RUBINSTEIN, A. (1998). *Modeling bounded rationality*, MIT press, Cambridge, USA.
- SAMET, D. (1990). Ignoring ignorance and agreeing to disagree, *Journal of Economic Theory* 52, 190–207.
- SANDHOLM, W. (2005). Excess payoff dynamics and other well-behaved evolutionary dynamics, *Journal of Economic Theory* 124, 149–170.
- (2006). *Population games and evolutionary dynamics*, MIT Press, forthcoming.
- SANDHOLM, W., DOKUMACI, E., LAHKAR, R. (2008). The projection dynamic and the replicator dynamic, *Games and Economic Behavior*, forthcoming.
- SEBENIUS, J., GEANAKOPOLOS, J. (1983). Don't bet on it: contingent agreements with asymmetric information, *Journal of the American Statistical Association* 78, 424–426.
- SIMON, H. (1955). A behavioral model of rational choice, *The Quarterly Journal of Economics* 69, 99–118.

- TAYLOR, P., JONKER, L. (1978). Evolutionarily stable strategies and game dynamics, *Mathematical Biosciences* 16, 76–83.
- VAN DAMME, E. (1991). *Stability and perfection of Nash equilibria*, 2nd ed. Springer, Berlin, Heidelberg, New York.
- VOORNEVELD, M. (2006). Probabilistic choice in games: Properties of Rosenthal’s t -solutions, *International Journal of Game Theory* 34, 105–121.
- WEIBULL, J. (1995). *Evolutionary Game Theory*, MIT Press, Cambridge, Massachusetts.