

April 22, 2012

“Evolutionary stability of Kantian optimization”

by

Philip A. Curry^a and John E. Roemer^b

Abstract. In Nash equilibrium, agents are autarchic in their optimization protocol, whereas in Kantian equilibrium, they optimize in an interdependent way. Typically, researchers into the evolution of *homo economicus* treat *preferences* as being determined by selective adaptation, but hold fixed the optimization protocol as autarchic. Here, we ask whether natural selection might choose the *optimizing protocol* to be either autarchic or interdependent. That is, will Kantian players, for whom the stable concept is Kantian equilibrium drive Nash players (for whom the stable concept is Nash equilibrium) to extinction, or otherwise? The answer depends upon whether players can signal their type to others.

Key Words: Kantian equilibrium, Nash equilibrium, evolutionary stable strategy, cooperation

JEL codes: C73, C62, D64

1. Introduction

There exists a literature in economics that considers some of the characteristics of agents, normally taken to be primitives, as being shaped by the forces of natural selection. This literature considers preferences¹ to be a heritable trait, and then examines how *homo*

^a Dept. of Economics, University of Waterloo

^b Depts. of Political Science and Economics, Yale University

¹ For example, economists have looked at the evolution of how individuals make decisions under uncertainty (see Robson [1996] and Curry [2001]), economic rationality (see Robson [2003a] and Gossner and Kuzmics [2009]), time preference (see Robson and Samuelson [2009] and Netzer [2009]) and altruism (see Bergstrom and Stark [1993] and Sethi and Somanathan [2001]). A more detailed summary of the literature can be found in Robson (2001, 2003b).

economicus would have evolved. The present paper examines the notion of Kantian optimization in strategic environments as developed in Roemer (1996, 2010, 2012) and asks whether there may be an evolutionary advantage to such thinking, particularly in relation to the autarchic optimization of Nash equilibrium. We find that Kantians hold an advantage over Nash players as long as they can identify whom they are playing against; absent this assumption, nature would select for Nash optimizers.

In particular, we examine the evolutionary dynamic in two games: the Prisoners' Dilemma and the Random Dictator Games. While players choose their strategies when playing these games, from an evolutionary standpoint, the meta-strategies are to behave as a Nash optimizer or a Kantian optimizer. We then look for Evolutionarily Stable Strategies (ESS's) as in Maynard Smith (1982, 1989). The determination of the type of player may arise from hard-wiring, or it may be a choice made by the individual based on past experiences and observations. In general, we are agnostic about the specifics concerning the determination of type, as the two forms of evolution operate in a similar fashion.

2. Kantian Equilibrium

The concept of Kantian equilibrium in games was introduced (see Roemer[1996, 2010, 2012]) to model the view that individuals in a society might apply the Kantian categorical imperative when deciding upon their actions. One of the simplest examples is that of a society of fishers, who own in common a lake, upon which they fish. Each fisher has preferences over fish consumed and labor expended in fishing, represented by a concave utility function $u^i(x, L)$. The total fish caught on the lake will be $G(\sum L^j)$, where G is a strictly concave production function of total labor expended. The catch going to fisher i will be proportional, statistically speaking, to the labor he expends, that is:

$$x^i(L^1, \dots, L^n) = \frac{L^i}{\sum L^j} G(\sum L^j). \quad (0.1)$$

Thus a game is defined whose strategies are fishing times and whose payoff functions are

$$V^i(L^1, \dots, L^n) = u^i(x^i(L^1, \dots, L^n), L^i). \quad (0.2)$$

It is well-known that, since G is strictly concave, the Nash equilibria of this game are Pareto inefficient; there is over-fishing, as each fisher ignores the congestion effect of his behavior on the productivity of labor on the lake, which is decreasing in total fishing time.

Suppose, however, that fishers reason in the following way. When facing a proposed labor vector $L = (L^1, \dots, L^n)$, fisher i thinks: “I have a right to increase (or decrease) my labor by a factor r if and only if I would be happy if all other fishers did so as well.” An equilibrium, with respect to this kind of behavior, is a vector L such that *nobody* would like *everybody* to alter his behavior by some factor $r \neq 1$. The formal definition is:

Definition 1. A vector $L = (L^1, \dots, L^n)$ is a *multiplicative Kantian equilibrium* of the fisher game V if and only if:

$$(\forall j)(\arg \max_{r \geq 0} V^j(rL^1, \dots, rL^n) = 1) \quad (0.3)$$

It is easy to prove that if a strictly positive vector L is a multiplicative Kantian equilibrium of the fisher game, then it is Pareto efficient (Roemer [2010]). In other words, Kantian behavior in the sense described solves the tragedy of the commons.

Elinor Ostrom (1990) has studied many small societies like the fisher economy, and discovers that they do *not* play the Nash equilibrium, but rather regulate fishing times in order to achieve Pareto efficient outcomes. These outcomes are ones in which each fisher keeps his catch. So they are efficient outcomes in which output (fish) is distributed in proportion to labor expended. It can be proved that *any such allocation is a multiplicative Kantian equilibrium* of the ‘game’ being played by fishers (see Roemer [2012]). The question arises: Have the fishers in these economies discovered the Kantian way of thinking? Does Kantian optimization play a role in maintaining these equilibria? We do not know. However, it is worth considering whether there may be an evolutionary explanation for the use of Kantian optimization.

3. Natural selection

In this section, we consider a large population of agents who are pair-wise matched at random each period to play a game; a fraction π of these agents optimize in the Nash manner, and $1 - \pi$ optimize in the Kantian manner. We examine the expected payoff to

each agent under two regimes. In the first regime, an agent's type is observable to his opponent. In the second, an agent's type is private knowledge. A discussion of these two regimes follows in section 4.

We consider payoffs to be proportional to evolutionary fitness, so that the type with the higher expected payoff will grow at a greater rate within the population. In this fashion, we can consider how a player optimizes and we examine whether an Evolutionarily Stable Strategy (ESS) exists, as in Maynard Smith (1982, 1989). We say that an ESS exists if either the expected payoff to Nash optimizers is greater when π is close to 1, or when the expected payoff to Kantian optimizers is greater when π is close to 0. In the former case, the population is almost all Nash optimizers and a mutant Kantian would do worse than the Nash players and so be driven to extinction. In the latter case, the population is almost all Kantian optimizers, and a Nash optimizer would do worse than the Kantians and so be driven out.

In the following two sections, we consider whether an ESS exists in two two-person games: the prisoners' dilemma, and the random dictator game.

3.1 The Prisoners' Dilemma Game

Consider the prisoners' dilemma (PD) game:

	Cooperate	Defect
Cooperate	(1,1)	(a,b)
Defect	(b,a)	(0,0)

where $b > 1$ and $a < 0$. Suppose that the Row player plays Cooperate with probability p , and the Column player Cooperates with probability q . Then the payoff functions of the two players are:

$$\begin{aligned} V^R(p,q) &= pq + p(1-q)a + (1-p)qb \\ V^C(p,q) &= pq + p(1-q)b + (1-p)qa \end{aligned} \tag{1.1}$$

A multiplicative Kantian equilibrium of this game is a strategy pair (p,q) such that

$$\arg \max_r V^R(rp, rq) = 1 = \arg \max_r V^C(rp, rq), \quad (1.2)$$

where r ranges over the interval $[0, \min[\frac{1}{p}, \frac{1}{q}]]$. We have:

Proposition 1

A. $(1,1)$ is a Kantian equilibrium of the PD game if and only if $(a+b) \leq 2$, and in this case, there is no other (non-trivial) Kantian equilibrium².

B. If $a+b > 2$, then the unique (non-trivial) Kantian equilibrium is given by

$$p^* = q^* = \frac{a+b}{2(a+b-1)} < 1.$$

In particular, $p^* > 1/2$.

Proof: See the appendix.

Proposition 1 gives us the payoffs when two Kantian optimizers play each other. We know that the only *Nash equilibrium* in the PD game is $p = q = 0$, since playing Defect is a strictly dominant strategy. We suppose that if a Kantian player knows that he faces a Nash player, he optimizes à la Nash: there is no point using the Kantian protocol as he knows the opponent will not use it. We first observe:

Proposition 2

Suppose each player can identify the type of his opponent. Then the Kantian players will drive the Nash players to extinction.

Proof:

1. If there is a match between two Kantian players, Proposition 1 tells us what the equilibrium is. The payoff to both players is positive regardless of the values of (a,b) .

² $(0,0)$ is the ‘trivial’ multiplicative Kantian equilibria of this game.

2. A Nash player always plays $p = 0$. So the payoff to two matched Nash players is zero.
3. Suppose a Nash player faces a Kantian player. Then they play $p = q = 0$ and their payoffs are each zero.
4. In sum, Nash players always receive zero payoffs in these matches, and Kantian players sometimes receive zero payoffs and sometimes positive payoffs. Consequently, on average, they are more fit than Nash players, and drive the Nash players to extinction. ν

We now study the case when matched players cannot identify the type of their opponent. In this case, a player assumes that with probability π his opponent is a Nash optimizer, and with probability $1 - \pi$ she is a Kantian optimizer. How should a Kantian player behave? Let the payoff function for the Kantian player be $P(p, q)$; the Nash player possesses a dominant strategy $q = 0$. If the Kantian player knows he is playing a Nash player, then he should simply play the Nash best response, since the player is not cooperating in the Kantian sense. However, suppose that if he plays another Kantian player, he anticipates will play the same strategy he is contemplating. Thus we claim the equilibrium play for a Kantian must be that probability p such that:

$$\arg \max_{0 \leq r \leq 1/x} \pi P(rp, 0) + (1 - \pi) P(rp, rp) = 1. \quad (2.4)$$

Expression (2.4) says that there is no strategy which does better for the Kantian which is both a unilateral best response to the Nash player and a Kantian (bilateral) best response to the Kantian player.

Proposition 3

If agents cannot identify the type of their matched opponent, then the Nash optimizers will drive to extinction the Kantian optimizers.

Proof:

1. Nash players have a dominant strategy: they always play $q = 0$.
2. If the first player is Kantian and the second is Nash, $(p, q) = (0, 0)$ is always an equilibrium. We shall assume that the Kantian player plays zero only if there is no other equilibrium.
3. A Kantian optimizer knows that if he faces a Nash player, then that player will play $q^N = 0$. Suppose a Kantian opponent plays q^K . The expected payoff to an agent playing p in this situation is, to paraphrase (2.4):

$$\pi p a + (1 - \pi)(p q^K + p(1 - q^K)a + (1 - p)q^K b). \quad (2.4a)$$

Suppose p is the equilibrium play of Kantian agent in this game. Then it must be the case that $p = q^K$, and so p must satisfy the following:

$$\arg \max_r \left(\pi r p a + (1 - \pi)(r^2 p^2 + r p(1 - r p)a + (1 - r p)r p b) \right) = 1. \quad (2.5)$$

The coefficient of r^2 in this expression is $(1 - \pi)p^2(1 - a - b)$. We have a several cases.

4. *Case 1.* If $a + b > 1$, then the maximand is a concave function of r , and achieves its maximum where its derivative w.r.t. r is zero. The statement that this derivative is zero at $r = 1$ is:

$$\pi p a + (1 - \pi)(2p^2 + p a - 2p^2 a + p b - 2p^2 b) = 0,$$

which solves to:

$$p^* = \frac{b(1 - \pi) + a}{2(1 - \pi)(-1 + a + b)}. \quad (2.6)$$

Indeed, the value p of (2.6) may not lie in the interval $[0, 1]$. So in this *Case* we have that the equilibrium play of a Kantian player is given by:

$$p = \begin{cases} 0, & \text{if } b(1 - \pi) + a \leq 0 \\ 1, & \text{if } b(1 - \pi) + a \geq 2(1 - \pi)(a + b - 1) \\ p^*, & \text{otherwise} \end{cases} \quad (2.7)$$

where the three cases correspond to the value of p^* being negative, greater than one, or in the interval $[0,1]$. Note that all three of these sub-cases can occur: the first occurs if π is close to 1, the second occurs if π is close to 0 and $a+b < 2$, and the third occurs if $a+b > 2$ and π is close to zero. We denote these three cases as 1a, 1b and 1c.

5. *Case 2.* $a+b \leq 1$. In this case, the maximand in (2.5) is a convex function of r , and so the maximum is achieved at an endpoint of feasible interval of factors r . One solution is $p=0$. If $p \neq 0$, then the interval of feasible values of r is $[0, 1/p]$. Since 1 must be the argmax, it follows that 1 is the right-hand endpoint of that feasible interval, and so $p=1$. It must also be the case that the value of the objective in (2.5) is weakly greater at 1 than at 0, which means:

$$\pi a + 1 - \pi \geq 0.$$

In sum, we have in this case:

- Case 2a: if $\pi a + 1 - \pi \geq 0$ then $p=1$ is an equilibrium
 - Case 2b: if $\pi a + 1 - \pi < 0$, then the only equilibrium is $p=0$.
6. We summarize these five cases as follows:

1a. $a+b > 1$ and $\pi > 1 + \frac{a}{b} \Rightarrow p^K = 0$

1b. $2 > a+b > 1$ and $\pi \leq \frac{2-(a+b)}{2-(a+b)-a} \Rightarrow p^K = 1$

1c. $a+b > 1$ and $0 \leq p^* \leq 1 \Rightarrow p^K = p^*$

2a. $a+b \leq 1$ and $\pi \leq \frac{1}{1-a} \Rightarrow p^K = 1$

2b. $a+b \leq 1$ and $\pi > \frac{1}{1-a} \Rightarrow p^K = 0$

(Regarding case 1b, note that if $a+b \leq 2$, note that the condition on π associated with the middle line of (2.7) is impossible.)

7. We must now compute the expected utility of Kantian players and Nash players in these five cases. Nash players always play zero. If a Nash player faces a Nash

player, the payoffs are zero. If a Nash player faces a Kantian player who plays p^K , his payoff is $p^K b$. So the expected utility of a Nash player is $(1-\pi)p^K b$.

We compute the expected utility of both players in the five cases:

Case 1a and 2b. Both Kantian and Nash players receive payoffs of zero.

Case 1b. The expected payoff of the Kantian player is greater than the Nash player if and only if:

$$\pi a + (1-\pi) > (1-\pi)b,$$

or $\pi > \frac{b-1}{a+b-1}$. In this case, this can happen if and only if there is a value of π such

that :

$$\frac{2-(a+b)}{2-(a+b)-a} \geq \pi > \frac{b-1}{a+b-1}.$$

This inequality reduces to the statement $2b(b-1) < a(1-a)$, an impossibility, since the l.h.s. is positive and the r.h.s. is negative. So in this case, the Nash players drive out the Kantian players.

Case 1c. In this case, the Kantian players have a higher payoff than the Nash players if and only if:

$$\pi p^* a + (1-\pi)p^* (p^* + (1-p^*)(a+b)) > (1-\pi)p^* b,$$

an inequality which reduces to the false statement : $a > (1-\pi)b$. So the Nash players drive out the Kantian players.

Case 2a. In this case, the expected payoff for a Kantian player is greater than that for a Nash player if and only if $\pi a + (1-\pi) > (1-\pi)b$, which is equivalent to the false statement $\pi(a+b-1) > b-1$. So the Nash players drive out the Kantian players.

8. We conclude that either Kantian players play $p^K = 0$, in which case their behavior is indistinguishable from that of Nash players, or they are driven to extinction. ν

3.2 The Random Dictator Game

In the random dictator game, there are two players, with von Neumann-Morgenstern utility functions $u(x) = \frac{x^a}{a}$, $v(x) = \frac{x^b}{b}$. The game has two stages. In stage 1, Nature chooses the dictator. In stage 2, the dictator keeps x and gives $1 - x$ to the opponent.

Thus, a strategy is a number $0 \leq x \leq 1$ for the a player and a strategy $0 \leq y \leq 1$ for the b player. The payoff functions are:

$$P^a(x, y) = \frac{1}{2}u(x) + \frac{1}{2}u(1 - y)$$

$$P^b(x, y) = \frac{1}{2}v(y) + \frac{1}{2}v(1 - x)$$

A *multiplicative Kantian equilibrium* is a pair (x, y) such that:

$$\frac{d}{dr} \Big|_{r=1} P^a(rx, ry) = 0 \text{ and } \frac{d}{dr} \Big|_{r=1} P^b(rx, ry) = 0.$$

In other words, neither player would like to multiply both strategies by any constant other than one. The following proposition proves that there is a unique Kantian equilibrium:

Proposition 4 *The unique multiplicative Kantian equilibrium of the random dictator game is $(\frac{1}{2}, \frac{1}{2})$.*

Proof: In the Appendix.

The unique Nash equilibrium to this game is to offer 0 when the dictator (i.e. keep 1 for oneself). As in the Prisoners' Dilemma game, it is assumed that a Kantian plays the Nash equilibrium when facing a Nash optimizer. Note again that we have:

Proposition 5 *In a population with Kantian and Nash agents, if players can identify their opponent's type, then Kantian players drive Nash players to extinction.*

Proof:

1. When Kantian players meet Kantian players, they each receive a payoff of $u(1/2)$, whereas when Nash players meet either Kantian or Nash players, both receive an expected payoff of $\frac{1}{2}u(1)$. By concavity, the Kantian players do better in expectation. v

We next study the equilibrium when each player knows she faces a Nash player with probability π and a Kantian player with probability $1 - \pi$.

Proposition 6. *In a population with a positive fraction of Nash players, where agents know the probability π but cannot identify their opponent's type, the Nash players drive the Kantian players to extinction.*

Proof:

1. We assume that both players have the same utility function, $u(x) = x^a / a$. The Nash player has a dominant strategy: to play $y = 1$.

2. Suppose it is a Kantian equilibrium for both Kantian players to play x in this situation. Then expression (2.4) becomes:

$$\arg \max_{0 \leq r \leq 1/x} \pi P^a(rx, 1) + (1 - \pi) P^a(rx, rx) = 1. \quad (3.3)$$

Equation (3.3) can be written:

$$\arg \max_{0 \leq r \leq 1/x} \frac{1}{2} u(rx) + \frac{1}{2} (\pi u(0) + (1 - \pi) u(1 - rx)) = 1. \quad (3.4)$$

We check whether there is a solution $0 < x < 1$. Since 1 is interior in the interval $[0, 1/x]$, and the expression in (3.4) is concave in r , we must have the FOC:

$$\frac{x}{2} u'(x) - \frac{1 - \pi}{2} x u'(1 - x) = 0, \quad (3.5)$$

or $u'(x) = (1 - \pi) u'(1 - x)$. For the utility function $u(x) = x^a / a$, we have:

$$x^* = \left(1 + (1 - \pi)^{\frac{1}{1-a}} \right)^{-1}.$$

Note that this checks: if $\pi = 0$, the equilibrium is that both players play $x = 1/2$.

3. Consequently, the expected utility of a Kantian player is $\frac{1}{2} u(x^*) + \frac{1 - \pi}{2} u(1 - x^*)$,

and the expected utility of a Nash player is $\frac{1}{2} u(1) + \frac{1 - \pi}{2} u(1 - x^*)$. Since $u(1) > u(x^*)$,

the Nash players drive the Kantian players to extinction. $\quad \vee$

4. Knowledge of Your Opponent

4.1 Secret Handshakes

The above section demonstrates that the success of Kantian optimizers hinges on their ability to determine whether the opponent is a fellow Kantian or a Nash player. In both games examined, Nash players do better if they are able to hide the fact that they are Nash optimizers, and so they have incentive to keep such information private.

This problem has been examined in the literature. Robson (1990) also considers agents who are pair-wise matched to play one-shot games with each other. He allows for agents to send a signal to each other, or alternatively, to have some characteristic that is observed before the game is played. Agents would then be able to condition their strategies on the signal received (or characteristic observed). In this context, we could consider Kantians to be sending such a signal, Robson's 'secret handshake'. If the two agents that have been paired together send the same signal, or possess the same characteristic, then each would recognize the other as a Kantian and behave accordingly.

Unfortunately, in Robson's analysis, this solution only works to select the Pareto dominant ESS in situations where there are multiple evolutionarily stable strategies. In both the Prisoners' Dilemma (a game examined by Robson) and the Random Dictator game, the unique ESS when the opponent's type is not known is to be a Nash optimizer, and so the secret handshake represents only a short-run solution to the problem. The reason is that a Nash player would like to discover the secret handshake so that he could convince Kantians that he was also Kantian, but then behave as a Nash player. In the long run, a mutant would indeed arise who would send the Kantian's signal but then play as a Nash optimizer.

That being said, Robson's secret handshake may be useful for coordination on the Kantian equilibrium. Without question, such signals cannot be useful for animals in the context of the two games we examine. However, we humans have constructed more complex environments for ourselves, and so may be able to take advantage of such signals in ways that other organisms cannot. In general, people rarely interact with others in such one-shot environments, but unless the repetition is infinite, what is true in a one-

shot game is also true if it is repeated. People also interact in many different ways, and have developed many ways of punishing others for inappropriate behavior.³

Without modeling such a complex interaction, in the games we consider it is worth noting that for Kantians to succeed, they should play the Kantian equilibrium strategy against fellow Kantians and the Nash equilibrium strategy against Nash optimizers. As such, one would not want to base any punishment purely on their play in the game, because Kantians *should* behave as Nash players if they think they are playing against a Nash player. Signals, however, would help others decide if a player should be punished for playing as a Nash player. If a player signaled that she were a Kantian and received a Kantian signal, then she should play the Kantian equilibrium. One would not be able to avoid punishment for playing Nash by claiming she thought she were playing against a Nash player in this context. Signals, then, could help establish the Kantian equilibrium when agents play these games within a larger context.

4.2 Group Selection

Group selection has not had a distinguished place in the evolutionary literature. Some early proponents of group selection were Sir Alexander Carr-Saunders (1922) and V. C. Wynne-Edwards (1962, 1963). However, the early general consensus was that group selection could not occur, although it could at times be indistinguishable from kin selection (see Maynard Smith [1964], Williams [1966], Dawkins [1976] and Grafen [1984]).

Yet the intuitive appeal of group selection is strong: groups that can avoid the deleterious effects of competition among themselves will do better than groups that succumb to inefficient outcomes induced by selfish behavior. For example, groups that solve the Prisoners' Dilemma will grow faster than groups that do not. Furthermore, humans are inherently social animals and in many ways the concept of the 'selfish gene' does not mesh with what we observe in our daily lives.

³ Lomas (2011) provides a survey of different ways that hunter gatherer societies attempt to deter undesirable behavior. Common methods include ostracism (exclusion for later interactions) and gossip (relaying of information about past behavior).

Bergstrom (2002) and Henrich (2004) provide overviews of how evolutionary theory has continually attempted to bring group selection into the fold. Some examinations of group selection could be applied to any species. First, it has been noted that ‘selfish genes’ often have incentive to do what’s best for the group (rather than what’s best for the individual) if the gene is prevalent within the group. In other words, when organisms tend to interact with genetic relatives, then there is little to distinguish between group selection and kin selection. Second, as alluded to above in the discussion of Robson’s secret handshake, if there exist mechanisms for punishment of anti-social behavior, such as defecting in the Prisoner’s Dilemma, then individuals’ incentives can be brought in line with the group’s (see Boyd and Richerson [1990,1992,2002] and Henrich and Boyd [2001]).

In addition to the above two concepts, Henrich (2004) considers that culture may play a role in facilitating group selection among humans. Culture can certainly influence who one interacts with as well as enable mechanisms for punishment of anti-social behavior. The effect of culture on the evolution of our species is not yet well understood, although it certainly appears to be important (see Boyd and Richerson [1985, 1996] and Henrich [2001]). This paper suggests that, in addition to kin selection and punishment, if culture can help reveal an agent’s true type then natural selection can favor what is good for the group rather than the individual.

5. Conclusion

We consider Kantian equilibrium in an evolutionary light. We find that the Kantian way of thinking can be evolutionarily advantageous *if* agents can observe the type of their opponent. This leads to the question of how agents could credibly signal their type, a problem that is not new to evolutionary biologists. Indeed, there is much in common with this problem and that of group selection. The results of this paper are somewhat novel in that we present a way for individual incentives to be aligned with social incentives while maintaining the assumption of purely self-regarding preferences.

To put the present paper in context, we can form a 2×2 taxonomy of models. With regard to *preferences*, individuals can be Selfish or Altruistic; with regard to *optimizing protocol*, individuals can be Nash or Kantian. The model of classical

economics and biology is {Selfish, Nash}. Much work in recent years has been done on the model {Altruistic, Nash}. In this paper, we have studied the model {Selfish, Kantian}. Why are the Kantian agents selfish in this paper? Because they are concerned only with their own utility, not the utility of others. Finally, we could study the model {Altruistic, Kantian}, in which agents care about others *and* optimize in the Kantian manner. This model is studied in Roemer (2012), although not in prisoners' dilemma or dictator games of this paper, but rather for large economies where allocation of the social product is the issue.

The present paper suggests that, while the literature has thus far looked to the evolution of altruistic preferences as a way of solving the group selection problem, some traction can be had by shifting focus to the optimizing protocol agents use.

Appendix

Proof to Proposition 1:

1. Note that V^C is monotone increasing in p and V^R is monotone increasing in q . Therefore, if (p, q) is a Kantian equilibrium where $p=0$ and $q \neq 0$ it must be that $q=1$ (because player 1 would advocate a change in scale factor of the efforts by $r = 1/q$) and if (p, q) is a Kantian equilibrium where $q=0$ and $p \neq 0$ then $p = 1$. It is easy to check that neither $(0,1)$ nor $(1,0)$ is Kantian. Therefore, the only non-trivial Kantian equilibria must be strictly positive.

Proof of Part A.

2. We write out the objective function for player R when he is choosing r :

$$V^R(rp, rq) = r^2 pq + rp(1-rq)a + (1-rp)rqb \quad (1.3)$$

We first show that $(1,1)$ is a Kantian equilibrium if and only if $a + b \leq 2$.

There are three cases to consider.

Case 1. $a + b > 1$. Then the objective function in (2.3) is concave in r , and so (2.3) is true if and only if the derivative of the objective is non-negative at 1, that is, when

$$a + b \leq 2.$$

Case 2. $a + b = 1$. Then the objective function is linear, and 1 is a maximum if and only if $a + b \geq 0$, which is true.

Case 3. $a + b < 1$. Then the objective function is convex in r , and 1 is a maximum if and only if its value at 1 is at least as large as its value at zero, that is:

$$(1 - a - b) + a + b \geq 0$$

which is surely true.

Putting these three cases together gives the result.

3. We next show that when $a + b \leq 2$, there is no other non-trivial Kantian equilibrium.

Case i). $2 \geq a + b \geq 1$ and $0 < p, q < 1$. The payoff functions are

$$V^R(rp, rq) = r^2 pq(1 - b - a) + r(pa + qb)$$

$$V^C(rp, rq) = r^2 pq(1 - a - b) + r(pb + qa)$$

which are concave functions of r when $a + b \geq 1$. If (p, q) is Kantian and p and q are less than one then the necessary and sufficient first-order conditions are:

$$2pq(1 - a - b) + pa + qb = 0$$

$$2pq(1 - a - b) + pb + qa = 0$$

which imply that $p = q$ and hence that $p = \frac{a + b}{2(a + b - 1)}$. But this value is at least one, a

contradiction.

4. Case ii). It is easy to check that there is no Kantian equilibrium when $a + b \geq 1$ and either p or q is one, unless both are one.
5. Case iii). $a + b < 1$. Since we know that p and q must both be positive in a (non-trivial) Kantian equilibrium (step 1), the functions $V^i(rp, rq)$ are strictly convex in r . Consequently, if neither p nor q is one, then the value $r = 1$ is in the interior of the interval of admissible scale factors, for both agents. But $V^i(rp, rq)$ must be maximized at a corner of the admissible interval of acceptable r 's. Therefore such a vector of probabilities cannot be Kantian. On the other hand, if $p = 1$, then player one will recommend a scale change of $r = 1/q$, since his payoff is increasing in the other's probability. Similarly if $q = 1$. Hence the only Kantian equilibrium is $(1, 1)$.

Proof of Part B.

6. Since $a + b > 2$, the functions $V^i(rp, rq)$ are strictly concave functions of r if $p \neq 0 \neq q$. Consequently, (p, q) is (non-trivial) Kantian only if either:

(i) at least one of $\{p, q\}$ is one, or

$$(ii) \left. \frac{dV^i}{dr}(rp, rq) \right|_{r=1} = 0.$$

We examine possibility (ii) first. This expands to the equations:

$$\begin{aligned} 2pq(1-a-b) + pa + qb &= 0 \\ 2pq(1-a-b) + pb + qa &= 0 \end{aligned}$$

which imply that $p = q$, which gives $p = \frac{a+b}{2(a+b-1)}$. As this number is positive and

less than one (since $a + b > 2$), this is indeed a Kantian equilibrium, as the statement in the proposition claims.

We investigate possibility (i). At most one of p or q can equal 1, by Part A of the proposition. Suppose $p = 1$. Then the first-order conditions for $(1, q)$ to be Kantian are:

$$\begin{aligned} 2q(1-a-b) + a + qb &\geq 0 \\ 2q(1-a-b) + b + qa &= 0 \end{aligned}$$

which imply that $a + bq \geq b + qa$ and therefore that $q \geq 1$, a contradiction. In like manner, it follows that $(p, 1)$ is not Kantian. $\quad \nu$

Proof to Proposition 4:

The FOC conditions defining K^\times equilibrium are:

$$\frac{x}{y} = \frac{u'(1-y)}{u'(x)}, \quad \frac{y}{x} = \frac{v'(1-x)}{v'(y)}. \quad (3.1)$$

Eliminating y from these equations, and using the explicit utility functions defined above, (3.1) can be written:

$$\frac{x^{\frac{1}{b}-a}}{(1-x)^{\frac{1}{b}-1}} = \left(1 - \frac{x^{\frac{1}{b}}}{(1-x)^{\frac{1}{b}-1}} \right)^{1-a} \quad (3.2)$$

We claim that $x = \frac{1}{2}$ is the unique solution of this equation. Note that the left-hand side

becomes $\left(\frac{1}{2}\right)^{1-a}$. So to demonstrate the claim, we need to show that $\frac{(1/2)^b}{(1/2)^{b-1}} = \frac{1}{2}$,

which is true. Furthermore, x is the unique solution, because the left-hand side is increasing in x , and the right-hand side is decreasing in x . ν

References

- Bergstrom, Theodore C. (2002) "Evolution of Social Behavior," *Journal of Economic Perspectives*, **16**, 67-88.
- Bergstrom, Theodore C. and Oded Stark (1993) "How Altruism Can Prevail in an Evolutionary Environment," *American Economic Review*, **83**, 149-155.
- Boyd, Robert and Peter Richerson (1985) *Culture and the Evolutionary Process*, Chicago: University of Chicago Press.
- (1990) "Group Selection Among Alternative Evolutionarily Stable Strategies," *Journal of Theoretical Biology*, **145**, 331-342.
- (1992) "Punishment Allows the Evolution of Cooperation (or Anything Else) in Sizable Groups," *Ethology and Sociobiology*, **13**, 171-195.
- (1996) "Why Culture is Common, but Cultural Evolution is Rare," *Proceedings of the British Academy*, **88**, 77-93.
- (2002) "Group Beneficial Norms Can Spread Rapidly in Structured Populations," *Journal of Theoretical Biology*, **215**, 287-296.
- Carr-Saunders, A.M. (1922) *The Population Problem: A Study in Human Evolution*, Oxford: Clarendon Press.
- Curry, Philip A. (2001) "Decision Making Under Uncertainty and the Evolution of Interdependent Preferences," *Journal of Economic Theory*, **98**, 357-369.
- Dawkins, Richard (1976) *The Selfish Gene*, Oxford: Oxford University Press.
- Gossner, Olivier and Christoph Kuzmics (2009) "Evolutionary Foundations of Rational Choice," working paper available at SSRN: <http://ssrn.com/abstract=1296656>
- Grafen, Alan (1984) "Natural Selection, Kin Selection and Group Selection," in *Behavioural Ecology, Second Edition*, J.R. Krebs and N.B. Davies, eds. London: Blackwell.
- Henrich, Joseph (2001) "Cultural Transmission and the Diffusion of Innovations: Adoption Dynamics Indicate that Biased Cultural Transmission is the Predominate Force in Behavioral Change and Much of Sociocultural Evolution," *American Anthropologist*, **103**, 992-1013.

----- (2004) “Cultural Group Selection, Coevolutionary Processes and Large-Scale Cooperation,” *Journal of Economic Behavior & Organization*, **53**, 3-35.

Henrich, Joseph and Robert Boyd (2001) “Why People Punish Defectors,” *Journal of Theoretical Biology*, **208**, 79-89.

Lomas, William (2009) “Conflict, Violence and Conflict Resolution in Hunting and Gathering Societies,” *Totem: The University of Western Ontario Journal of Anthropology*, **17**, Article 13.

Maynard Smith, John (1964) “Group Selection and Kin Selection,” *Nature*, **201**, 1145-1147.

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

----- (1989) *Evolutionary Genetics*, Oxford: Oxford University Press.

Netzer, Nick (2009) “Evolution of Time Preferences and Attitudes Towards Risk,” *American Economic Review*, **99**, 937-955.

Ostrom, E. 1990. *Governing the commons: The evolution of institutions for collective action*, Cambridge UK: Cambridge University Press

Robson, Arthur J. (1990) “Efficiency in Evolutionary Games: Darwin, Nash and the Secret Handshake,” *Journal of Theoretical Biology*, **144**, 379-396.

----- (1996) “A Biological Basis for Expected and Non-Expected Utility,” *Journal of Economic Theory*, **68**, 190-207.

----- (2001) “The Biological Basis of Economic Behavior,” *Journal of Economic Literature*, **39**, 11-33.

----- (2003a) “The Evolution of Rationality and the Red Queen” *Journal of Economic Theory*, **111**, 1-22.

----- (2003b) “Evolution and Human Nature,” *Journal of Economic Perspectives*, **16**, 89-106.

Robson, Arthur J. and Larry Samuelson (2009) “The Evolution of Time Preference with Aggregate Uncertainty,” *American Economic Review*, **99**, 1925-1953.

Roemer, J.E. 1996. *Theories of distributive justice*, Cambridge MA: Harvard University Press

----- 2010. “Kantian equilibrium,” *Scandinavian J. Econ.* **112**, 1-24

----- 2012. "Kantian optimization, social ethos, and Pareto efficiency" ,Cowles Foundation Discussion Paper No. 1854, Yale University

Sethi, Rajiv and E. Somanathan (2001) "Preference Evolution and Reciprocity," *Journal of Economic Theory*, **97**, 273-297.

Williams, George C. (1966) *Adaptation and Natural Selection: A Critique of Some Current Evolutionary Thought*, Princeton: Princeton University Press.

Wynne-Edwards, V.C. (1962) *Animal Dispersion in Relation to Social Behaviour*, London: Oliver and Boyd.

----- (1963) "Intergroup Selection in the Evolution of Social Systems," *Nature*, **200**, 623-626.