

Samaritan vs Rotten Kid: Another Look

Bouwe R. Dijkstra*

University of Nottingham

January 2002

Abstract. We set up a two-stage game with sequential moves by one altruistic agent and n selfish agents. The rotten kid theorem states that the altruist can only reach her first best when the selfish agents move before the altruist. The Samaritan's dilemma, on the other hand, states that the altruist can only reach her first best when she moves before the selfish agents. We find that in general, the altruist can reach her first best when she moves first, if and only if a selfish agent's action marginally only affects his own payoff. The altruist can reach her first best when she moves last if and only if there is just one commodity involved. When the altruist cannot reach her first best when she moves last, the outcome is not Pareto efficient either.

JEL Classification: D64

Key words: Altruism, rotten kid theorem, Samaritan's dilemma

Lear: I gave you all—

Regan: And in good time you gave it.

(King Lear, Act II, Scene IV)

Correspondence: Bouwe R. Dijkstra, School of Economics, University of Nottingham, Nottingham NG7 2RD, UK. Tel: +44 115 8467205, Fax: +49 115 9514159, email: bouwe.dijkstra@nottingham.ac.uk

*I thank Arye Hillman, Gordon Tullock and Heinrich Ursprung for pointing me toward the subject of this paper. I thank Richard Cornes, Hans Gersbach, Johan Lagerlöf, Andreas Lange, Markus Lehmann, Andries Nentjes, Till Requate, Bert Schoonbeek and Perry Shapiro for valuable comments. Special thanks are due to Ted Bergstrom for essential clarifications.

1 Introduction

However much we care about other people, we do not wish to invite them to take advantage of our charity. The economic theory of altruism offers two conflicting pieces of strategic advice: the rotten kid theorem (Becker [1] [2]) and the Samaritan's dilemma (Buchanan [7]). In a single-round model with sequential moves by an altruistic agent (the Samaritan or the parent) and a selfish agent (the parasite or the kid), the contradiction between the two can be stated as follows.

The rotten kid theorem states that the parent can only reach her first best when she moves after the kid. The intuition is that the kid will only act unselfishly if the parent can reward him afterward. The Samaritan's dilemma, on the other hand, states that the Samaritan can only reach his first best when he moves before the parasite. Here, the intuition is that only when the Samaritan moves first will his actions be immune to manipulation by the parasite.

In this paper, we shall identify the restrictions on the agents' payoff functions for either result to hold. For the altruist to reach her first best when she moves first, a selfish agent's actions should only affect his own payoff on the margin. Then there are no externalities to his actions. For the altruist to reach her first best when she moves last, there should only be one commodity, which we might call income. Then a selfish agent cannot manipulate the altruist's trade-off between her own and the selfish agents' payoffs. The selfish agents will maximize aggregate income. They benefit from this themselves, because their payoffs are normal goods to the altruist.

As we interpret Samaritan's dilemma and rotten kid theorem, they have a positive as well as a negative side. The positive side is that the altruist can reach her first best under one sequence of moves. The negative side is that she cannot reach her first best under the other sequence. Many authors have used the terms Samaritan's dilemma and rotten kid theorem in the positive sense only. We shall refer to these versions as the positive Samaritan's dilemma and the positive rotten kid theorem.

Our result for the positive Samaritan's dilemma is new. For the positive rotten kid theorem, Bergstrom [3] has performed a similar analysis. His model can be seen as a

special version of our more general setup. Whereas we do not restrict the nature of the altruist's actions, Bergstrom [3] assumes she distributes a certain amount of money among the selfish agents. Removing this restriction results in a slightly more general condition for the positive rotten kid theorem.

The focus of this paper is on the simple one-shot game with complete information with which the theory started 25 year ago. Since then, more complex games between altruists and selfish agents have been studied.¹ It would be worthwhile to expand the general analysis to encompass multi-period models and asymmetric or incomplete information.

However peripheral to economics the study of altruism may seem, there is in fact an application that takes us to the very heart of the discipline (Munger [22]). Regarding the welfare-maximizing government as an altruist and the private agents as selfish agents, we have a framework for a policy game. This framework allows us to study how the government can shape incentives such that private actions maximize social welfare.

The rest of this paper is organized as follows. In Section 2, we introduce and discuss the Samaritan's dilemma and the rotten kid theorem in simple two-agent setups where they are known to hold. In Section 3, we set up a single-round game with n selfish agents, deriving the conditions for the Samaritan's dilemma and the rotten kid theorem to hold. In Section 4, we discuss Samaritan's dilemma and rotten kid theorem in terms of Pareto efficiency. In Section 5, we discuss Bergstrom's [3] game as well as Bergstrom's [3] own and Cornes and Silva's [11] conditions for the rotten kid theorem. We conclude with Section 6.

2 Introductory examples

2.1 Samaritan's dilemma

The Samaritan's dilemma is due to Buchanan [7] who discusses a game between an altruistic Samaritan and a selfish parasite.² He shows that the Samaritan can reach his first

¹Bruce and Waldman [5] [6] and Lindbeck and Weibull [19] have analyzed two-period lifetime models. Chami [9] [10] and Lagerlöf [18] assume asymmetric information. Coate [8], Lord and Raganzas [20] and Wigger [26] include uncertainty.

²Buchanan [7] distinguishes between the active and the passive Samaritan's dilemma. We shall only discuss the passive Samaritan's dilemma here. The passive Samaritan's preferences are reconcilable with a payoff function that only depends on his donation. The active Samaritan's payoff, on the other hand,

best when he moves before the parasite, but not when he moves after the parasite. In this subsection, we shall present a continuous version of the game.³

The Samaritan maximizes his objective function $W(U_0, U_1)$, increasing in his own payoff U_0 and the parasite's payoff U_1 : $W_k \equiv \partial W / \partial U_k > 0, k = 0, 1$. The parasite maximizes his own payoff U_1 . The Samaritan's own payoff U_0 only depends on his donation y to the parasite, so that we can simply set $U_0 = -y$. The parasite's payoff depends on his work effort x and on the Samaritan's donation y . The parasite's payoff function $U_1(y, x)$ has the following properties:

- $\partial U_1 / \partial y > 0, \partial^2 U_1 / \partial y^2 \leq 0$. The parasite's marginal payoff of money is positive and decreasing.
- $\partial U_1 / \partial x > [<]0$ for $x < [>]x^*(y), x^*(y) > 0; \partial^2 U_1 / \partial x^2 \leq 0$. Given the Samaritan's donation y , there is an optimal work effort $x^*(y)$ for the parasite, where the marginal payoff of extra money earned equals the marginal payoff of leisure.
- $\partial^2 U_1 / \partial y \partial x < 0$. An increase in the parasite's effort decreases his marginal payoff of money. This is because the parasite earns more money when he works harder and his marginal payoff of money is decreasing.

The first order conditions for the Samaritan's first best are, with respect to y and x , respectively:

$$W_0 = W_1 \frac{\partial U_1}{\partial y} \tag{1}$$

$$\frac{\partial U_1}{\partial x} = 0 \tag{2}$$

We shall now see that the Samaritan can always reach his first best when he moves first, but he can never reach his first best when he moves last.

must also depend on the parasite's action. This follows from the fact that, given that the Samaritan donates, the active Samaritan prefers the parasite to go to work although the parasite prefers to stay in bed. Schmidtchen [24] provides an analysis of the active Samaritan's dilemma.

³Jürges [16] also analyzes this game. Bergstrom ([3], 1140-1) analyzes a similar game, where a parent distributes money after his "lazy rotten kids" have set their work efforts. Neither Bergstrom [3] nor Jürges [16] identify the game with the Samaritan's dilemma.

When the Samaritan moves first, the parasite sets x in stage two to maximize his own payoff:

$$\frac{\partial U_1}{\partial x} = 0$$

This condition is identical to the first order condition (2) for the Samaritan's first best with respect to x . Thus, in stage one, the Samaritan can set y according to his first best condition (1). This means that the Samaritan can always reach his first best when he moves first.

The intuition is that the parasite sets the work effort that maximizes his own payoff, taking the Samaritan's donation as given. Since the parasite's work effort only affects his own payoff, the parasite takes the full effect of his decision into account. There is no externality, and the Samaritan's first best is implemented.

When the parasite moves first, the Samaritan sets y according to (1) in stage two. In stage one, the parasite sets the x that maximizes his own payoff, taking into account that his choice of x affects the Samaritan's choice of y in stage two:

$$\frac{dU_1}{dx} \equiv \frac{\partial U_1}{\partial x} + \frac{\partial U_1}{\partial y} \frac{dy}{dx} = 0$$

This only corresponds to the Samaritan's first order condition (2) for x when $dy/dx = 0$, i.e. the donation reaches its maximum, in the optimum. In order to find the expression for dy/dx in the optimum, we totally differentiate the Samaritan's first order condition for y (1) with respect to x and substitute (2):

$$\frac{dy}{dx} = \frac{W_1 \left(\frac{\partial^2 U_1}{\partial y \partial x} \right)}{-W_{00} + (W_{10} + W_{01}) \frac{\partial U_1}{\partial y} - W_{11} \left(\frac{\partial U_1}{\partial y} \right)^2 - W_1 \frac{\partial^2 U_1}{\partial y^2}} < 0 \quad (3)$$

The numerator in (3) is negative, because $W_1 > 0$ and $\partial^2 U_1 / \partial y \partial x < 0$. The denominator is positive, because this is the second order condition $\partial^2 W / \partial y^2 < 0$.

Thus, the parasite gets more money from the Samaritan, the less he works. As a result, the parasite will work less than the Samaritan would like him to. The Samaritan cannot reach his first best when he moves after the parasite. Intuitively, the less money the parasite earns, the needier he is and the more money he will get from the Samaritan.

When the parasite moves first, he can extort more money from the Samaritan by working less. We can also say that the parasite gets the Samaritan to buy more of his payoff U_1 by lowering its price.

2.2 Rotten kid theorem

In order to introduce the rotten kid theorem, we analyze the simple game discussed by Becker [1] [2] and commented upon by Hirshleifer [15]. The game is between an altruistic parent and a selfish kid. The kid can undertake an action that affects his own as well as the parent's income. The parent can give money to the kid. We shall see that in general, the parent cannot reach her first best when she moves first, but she can always reach her first best when she moves after the kid.

In fact, Becker [1] [2] himself does not discuss the order of moves. Citing Shakespeare's *King Lear*, Hirshleifer [15] was the first to point out that the parent's first best is implemented only when the kid moves first.⁴

Denote the kid's action by x and the parent's transfer by y . Since the only commodity involved is income, we can equate the parent's and kid's payoffs, U_0 and U_1 respectively, with income and write them in the additively separable form:

$$U_0 = -y + b_0(x) \quad U_1 = y + b_1(x) \quad (4)$$

Here, $b_k(x)$, $k = 0, 1$, is the effect of the kid's action on the income of the parent and the kid, respectively.

The selfish kid maximizes his own payoff U_1 . The parent maximizes her objective function $W(U_0, U_1)$ with $W_k \equiv \partial W / \partial U_k > 0$, $k = 0, 1$.

The first order conditions for the parent's first best are, with respect to y and x

⁴Pollak [23] offers an alternative qualification: The parent can reach her first best only if she makes a take-it-or-leave-it offer to the kid. The offer specifies the kid's action and the parent's transfer. Cox [12] elaborates on this point. He argues that the parent can only reach her first best if the kid is better off accepting the offer to implement the first best than rejecting it. Cox [12] calls this "altruism". If the kid's participation constraint is binding, the parent will offer a different contract which gives the kid his reservation payoff. Cox calls this "exchange". In our model, we assume that the selfish agent's participation constraint never binds.

respectively:

$$W_0 = W_1 \tag{5}$$

$$W_0 b'_0 + W_1 b'_1 = 0 \tag{6}$$

Substituting (5) into (6):

$$b'_0 + b'_1 = 0 \tag{7}$$

This implies that in the parent's first best, family income $U_0 + U_1 = b_0 + b_1$ is maximized.

When the parent moves before the kid, the kid will set $b'_1 = 0$. In general, this does not correspond to the parent's first order condition (7). When the kid moves last, he will maximize his own income instead of family income.

Now we shall see what happens when the kid moves before the parent. In stage two, the parent will set the transfer y that maximizes W , according to (5). In stage one, the kid sets the x that maximizes his income, taking into account that his action affects the parent's transfer:

$$\frac{dU_1}{dx} \equiv \frac{dy}{dx} + b'_1 = 0 \tag{8}$$

The value of dy/dx follows from the total differentiation of the parent's first order condition (5) with respect to x :

$$(W_{00} - W_{10}) \left(-\frac{dy}{dx} + b'_0 \right) = (W_{11} - W_{01}) \left(\frac{dy}{dx} + b'_1 \right) \tag{9}$$

By the kid's first order condition (8), the second term between brackets on the RHS of (9) is zero. Thus, the second term between brackets on the LHS of (9) must be zero:

$$\frac{dy}{dx} = b'_0$$

Substituting this into the kid's first order condition (8), we see that it is equivalent to the parent's first best condition (7): the kid effectively maximizes family income.

Thus, the parent always reaches her first best when she moves after the kid. As Bernheim et al. [4] and Bergstrom [3] already noted, this result follows from the assumption that there is only one commodity, namely income. The intuition, due to Bergstrom [3],

is that when there is only one commodity, say income, we can identify payoff with income. The kid cannot manipulate the price of his income in terms of the parent's income, because it is always unity. Then the parent and the kid agree that it is a good thing to maximize aggregate income. It is clear that the parent will want to maximize family income. However, as Becker [1] already notes, the kid will only want to maximize family income if he benefits from that himself, i.e. if his payoff is a normal good to the parent.⁵

3 A general analysis

3.1 The model

In this section, we analyze a model with one altruistic agent and n selfish agents. We shall see under which conditions the Samaritan's dilemma and the rotten kid theorem hold.

There are $n + 1$ agents, indexed by $k = 0, \dots, n$. Agent 0 is the altruist and agents i , $i = 1, \dots, n$, are the selfish agents. Agent i controls the variable x_i . Agent 0 can contribute to each agent i 's payoff U_i . This contribution is denoted by y_i . Therefore, $\partial U_i / \partial y_i > 0$ for $i = 1, \dots, n$, and $\partial U_i / \partial y_j = 0$ for all $j = 1, \dots, n$, $j \neq i$, by definition.

There will be an upper and a lower bound to $\mathbf{y} = (y_1, \dots, y_n)$. The lower bound is $\mathbf{y} = 0$: agent 0 can only give to the other agents, she cannot improve her own payoff at the expense of the others. The upper bound follows from the restriction that agent 0 only has a limited amount of time, money, or whatever the nature of \mathbf{y} , to give to the others. The exact formulation of the upper bound depends on the nature of \mathbf{y} . We shall assume that neither the upper nor the lower bound are binding constraints on the equilibria.

Agent 0's payoff has the form $U_0(\mathbf{y}, \mathbf{x})$, which is continuous and twice differentiable, with $\mathbf{x} = (x_1, \dots, x_n)$. Agent i 's payoff has the form $U_i(y_i, \mathbf{x})$, which is continuous and twice differentiable with $\partial^2 U_i / \partial x_i^2 \leq 0$. Each agent i , $i = 1, \dots, n$, maximizes his own payoff. Agent 0, however, does not only care about her own payoff, but also about the payoffs of all other n agents. Her objective function is $W(\mathbf{U})$, continuous and twice differentiable with $\mathbf{U} \equiv (U_0, \dots, U_n)$, $W_k \equiv \partial W / \partial U_k > 0$, $k = 0, \dots, n$.

Let us now determine the first-best outcome for agent 0. We assume that the first best is characterized by an interior solution. Thus, W should be concave in (\mathbf{y}, \mathbf{x}) . Dif-

⁵A formal proof of this point in the general setup of Section 3 is available from the author.

differentiating $W(\mathbf{U})$ with respect to y_i , $i = 1, \dots, n$, we find the following n first order conditions:

$$W_0 \frac{\partial U_0}{\partial y_i} + W_i \frac{\partial U_i}{\partial y_i} = 0 \quad (10)$$

Note that since $W_0, W_i > 0$ and $\partial U_i / \partial y_i > 0$, we must have $\partial U_0 / \partial y_i < 0$ in the optimum. Differentiating $W(\mathbf{U})$ with respect to x_i , $i = 1, \dots, n$, we find the following n first order conditions:

$$\sum_{k=0}^n W_k \frac{\partial U_k}{\partial x_i} = 0 \quad (11)$$

Whatever agent 0's precise preferences, her first best will always be on the payoff possibility frontier PPF . Every element \mathbf{U}^* of the PPF is defined as:

$$U_i^*(\mathbf{U}_{-i}^*) \equiv \max_{\mathbf{x}, \mathbf{y}} U_i \text{ s.t. } \mathbf{U}_{-i} = \mathbf{U}_{-i}^*, \quad i = 1, \dots, n \quad (12)$$

where $\mathbf{U}_{-i} \equiv (U_0, \dots, U_{i-1}, U_{i+1}, \dots, U_n)$. Let \mathbf{x}^* be an \mathbf{x} vector that is associated with a \mathbf{U}^* , and \mathbf{X}^* the set of all \mathbf{x}^* :

$$\begin{aligned} \mathbf{x}^*(\mathbf{U}^*) &= \arg \max_{\mathbf{x}} U_i \text{ s.t. } \mathbf{U}_{-i} = \mathbf{U}_{-i}^* \\ \mathbf{X}^* &\equiv \{\mathbf{x}^*(\mathbf{U}^*)\} \end{aligned} \quad (13)$$

In the following, we shall study the effect of sequential moves. The agents i , $i = 1, \dots, n$, will always move simultaneously. In subsection 3.2, we see what happens when agent 0 moves before agents i . In subsection 3.3, we analyze the case where the agents i move before agent 0. We will derive the conditions for these sequences of moves to result in agent 0's first best for all $W(\mathbf{U})$. The conditions will thus be on the payoff functions \mathbf{U} . We are looking for the necessary and sufficient local restrictions on \mathbf{U} under which the first order conditions of the subgame perfect equilibrium are equal to the first order conditions (10) and (11) of agent 0's first best. We shall assume that the second order conditions, which involve a combination of restrictions on $W(\mathbf{U})$ and \mathbf{U} , are satisfied.

In our interpretation of the Samaritan's dilemma and the rotten kid theorem, they do not only have a positive side to them (agent 0 can reach her first best under one sequence of moves), but also a negative side: Agent 0 cannot reach her first best under the other sequence. In subsection 3.4, we give the formal definitions and state the conditions for the Samaritan's dilemma to apply and for the rotten kid theorem to hold.

3.2 Agent 0 moves first

In this subsection, we derive the equilibrium for the game where agent 0 moves before agents i , and we see when this equilibrium corresponds to the first best for agent 0. Thus, we shall derive the condition for the positive Samaritan's dilemma to hold:

Definition 1 *The positive Samaritan's dilemma states that agent 0 can reach her first best when she moves in stage one and agents i , $i = 1, \dots, n$, move in stage two.*

We assume an interior solution. The game is solved by backwards induction. In stage two, each agent i , $i = 1, \dots, n$, sets the x_i that maximizes his own payoff, taking y_i and all other x_l , $l = 1, \dots, i - 1, i + 1, \dots, n$, as given:

$$\frac{\partial U_i}{\partial x_i} = 0 \quad (14)$$

In stage one, agent 0 sets the y_i that maximize her objective function $W(\mathbf{U})$, taking into account that the agents i 's choices of x_i depend upon her choice of y_i :

$$W_0 \frac{\partial U_0}{\partial y_i} + W_i \frac{\partial U_i}{\partial y_i} + \sum_{k=0}^n W_k \frac{\partial U_k}{\partial x_i} \frac{dx_i}{dy_i} = 0$$

Substituting (14) and differentiating the i th condition (14) totally with respect to y_i , this can be rewritten as:

$$W_0 \frac{\partial U_0}{\partial y_i} + W_i \frac{\partial U_i}{\partial y_i} - \sum_{\substack{l=0 \\ l \neq i}}^n W_l \frac{\partial U_l}{\partial x_i} \frac{\partial^2 U_i / \partial y_i \partial x_i}{\partial^2 U_i / \partial x_i^2} = 0 \quad (15)$$

In general, the outcome will not be agent 0's first best. We shall now see under which condition agent 0 can reach her first best when she moves first.⁶

Condition 1 *For all $\mathbf{x} \in \mathbf{X}^*$, all $j = 1, \dots, n$ and all $l = 0, \dots, n$, $l \neq j$:*

$$\frac{\partial U_l}{\partial x_j} = 0$$

Proposition 1 *Given that all agents' second order conditions are satisfied, the positive Samaritan's dilemma holds for all $W(\mathbf{U})$ if and only if Condition 1 holds.*

⁶All proofs are in the Appendix.

The intuition behind the result is straightforward. When selfish agent i moves last, he does not take into account the effect of his action on any of the other agents' payoffs. In general, this can only result in the first best for agent 0 if agent i 's action does not affect any other agent's payoff,⁷ at least not on the margin. Then agent i takes the full effect of his actions into account. There is no externality, and agent 0's first best is implemented.

3.3 Agents i move first

In this subsection, we derive the equilibrium for the game where agents i move before agent 0, and we see when this equilibrium corresponds to the first best for agent 0. Thus, we shall derive the conditions for the positive rotten kid theorem to hold:

Definition 2 *The positive rotten kid theorem states that agent 0 can reach her first best when agents i , $i = 1, \dots, n$, move in stage one and agent 0 moves in stage two.*

We solve the game by backwards induction, assuming an interior solution. In stage two, agent 0 sets the y_j that maximize her objective function $W(\mathbf{U})$, taking all x_j , $j = 1, \dots, n$, as given:

$$W_0 \frac{\partial U_0}{\partial y_j} + W_j \frac{\partial U_j}{\partial y_j} = 0 \quad (16)$$

Of course these conditions are identical to the first order conditions (10) for agent 0's first best with respect to y_j .

In stage one, each agent i , $i = 1, \dots, n$, sets the x_i that maximizes his own payoff, taking the x_l , $l = 1, \dots, i-1, i+1, \dots, n$, from the other $n-1$ agents moving in stage one as given, but realizing that his choice of x_i affects agent 0's choice of y_i in stage two:

$$\frac{dU_i}{dx_i} \equiv \frac{\partial U_i}{\partial y_i} \frac{dy_i}{dx_i} + \frac{\partial U_i}{\partial x_i} = 0 \quad (17)$$

The values for dy_j/dx_i , $j = 1, \dots, n$, follow from the total differentiation of the n conditions (16) with respect to x_i . We shall write the total differential of (16) with respect to x_i in a compact manner that will prove useful later:

$$\frac{\partial U_0}{\partial y_j} \sum_{k=0}^n W_{0k} \frac{dU_k}{dx_i} + W_0 \frac{d(\partial U_0 / \partial y_j)}{dx_i} + \frac{\partial U_j}{\partial y_j} \sum_{k=0}^n W_{jk} \frac{dU_k}{dx_i} + W_j \frac{d(\partial U_j / \partial y_j)}{dx_i} = 0 \quad (18)$$

⁷This is the condition we already encountered in subsection 2.1.

In general, the equilibrium conditions (17) and (18) for $x_i, i = 1, \dots, n$, are not identical to the corresponding first order conditions (11) for agent 0's first best. We shall now see when they are. Let us first state an intermediate result.

Condition 2 *Consider a marginal change in x_i after which \mathbf{y} is adjusted optimally according to (16). Define*

$$\frac{dU_0}{dx_i} \equiv \frac{\partial U_0}{\partial x_i} + \sum_{j=1}^n \frac{\partial U_0}{\partial y_j} \frac{dy_j}{dx_i} \quad (19)$$

$$\frac{dU_j}{dx_i} \equiv \frac{\partial U_j}{\partial x_i} + \frac{\partial U_j}{\partial y_j} \frac{dy_j}{dx_i} \quad (20)$$

for all $i, j = 1, \dots, n$. Then $dU_l/dx_i = 0$ when $dU_i/dx_i = 0$ for all $\mathbf{x} \in \mathbf{X}^*$ and for all $l = 0, \dots, n, i = 1, \dots, n, l \neq i$.

Lemma 1 *Given that the second order conditions are satisfied, the positive rotten kid theorem holds for all $W(\mathbf{U})$ if and only if Condition 2 holds.*

Note the analogy with Proposition 1 from subsection 3.2. When agents i move last, they set $\partial U_i/\partial x_i = 0$. This will result in agent 0's first best for all $W(\mathbf{U})$ if and only if $\partial U_l/\partial x_i = 0$ for all $l \neq i$. When agents i move first, they set $dU_i/dx_i = 0$. This will result in agent 0's first best for all $W(\mathbf{U})$ if and only if $dU_l/dx_i = 0$ for all $l \neq i$.

Condition 2 is not a condition on the payoff functions yet. It requires agent 0 to set all $dU_l/dx_i = 0$ and thus depends on agent 0's behaviour.

Condition 3 \mathbf{X}^* *consists of a single vector $\tilde{\mathbf{x}}$.*

Proposition 2 *Given that the second order conditions are satisfied, the positive rotten kid theorem holds for all $W(\mathbf{U})$ if and only if Condition 3 holds.*

We can say that the price of an agent i 's payoff along the *PPF* should be constant and beyond manipulation by agent i . Then we can aggregate all payoffs using these prices and refer to aggregate payoff as income $I(\mathbf{x})$, as defined in (??). The agents i will maximize income and agent 0 will redistribute income. In the terminology of Monderer and Shapley [21], Condition 3 turns the game into a potential game, where all agents $i = 1, \dots, n$ maximize the ordinal potential function $I(\mathbf{x})$.

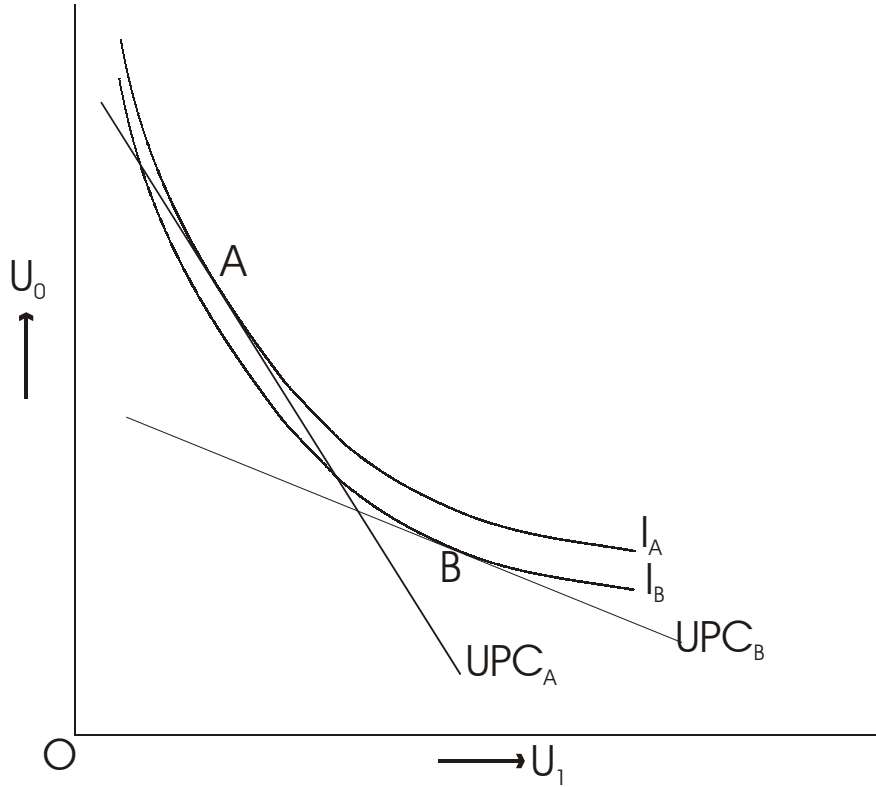


Figure 1: Intersecting Utility Possibility Curves

Figure 1 illustrates what goes wrong when a selfish agent can influence the price of his payoff or equivalently, when the *PPF* consists of multiple *UPCs*. Point *A* is agent 0's first best. It is reached when the single selfish agent 1 selects the action x^A that implements UPC_A . Let UPC_A be the anchor according to Definition ??, so that it is a straight line.

Now suppose agent 1 can decrease the price of his own payoff, either by increasing or decreasing his x . For instance, when agent 1 chooses x^B , the resulting UPC_B is flatter than UPC_A , lies everywhere below agent 0's first-best indifference curve I_A and intersects UPC_A so that the *PPF* does not consist of UPC_A alone. At the point where UPC_B comes closest to I_A , U_1 is higher and U_0 is lower than in *A*, because UPC_B is flatter than UPC_A . Then U_1 will also be higher in point *B*, where agent 0's indifference curve I_B is tangent to UPC_B . Agent 1 will prefer implementing UPC_B to UPC_A , because U_1 is higher in point *B* on UPC_B than in point *A* on UPC_A .

3.4 Conditions for Samaritan's dilemma and rotten kid theorem

There are two alternative definitions for the Samaritan's dilemma and the rotten kid theorem, applied to agent 0's first best. The positive definition of the Samaritan's dilemma (the rotten kid theorem) is: agent 0 can reach her first best when she moves first (last). In subsections 3.2 and 3.3, we have already defined the positive versions of Samaritan's dilemma and rotten kid theorem and derived the conditions for them to hold.

The second definition of the Samaritan's dilemma (the rotten kid theorem) also includes a negative side: agent 0 can reach her first best when she moves first (last), but not when she moves last (first). This is the definition we adhere to in this paper. We shall now formally define the Samaritan's dilemma and the rotten kid theorem:

Definition 3 *The Samaritan's dilemma states that agent 0 can reach her first best when she moves in stage one and agents i , $i = 1, \dots, n$, move in stage two, but not when agents i move in stage one and agent 0 moves in stage two.*

Definition 4 *The rotten kid theorem states that agent 0 can reach her first best when agents i move in stage one and agent 0 moves in stage two, but not when agent 0 moves in stage one and agents i , $i = 1, \dots, n$, move in stage two.*

For the Samaritan's dilemma, we can simply take our Conditions 1 and 3:

Proposition 3 *Given that all agents' second order conditions are satisfied, the Samaritan's dilemma holds for all $W(\mathbf{U})$ if and only if Condition 1 holds and Condition 3 does not hold.*

For the rotten kid theorem, the analysis is somewhat more complicated. Agent 0 can reach her first best when she moves first for any $W_k > 0$ if all $\partial U_l / \partial x_j = 0$, $j = 1, \dots, n$, $l = 0, \dots, n$, $l \neq j$ (Condition 1). But under Condition 3, which ensures that agent 0 can reach her first best when she moves last, substituting (??) into (10) reveals that all W_k are equal in the first best. Then agent 0 can reach her first best when she moves first if all $\sum_l \partial U_l / \partial x_j = 0$. We don't need $\partial U_l / \partial x_j = 0$ for all l , as long as the sum is zero.⁸

⁸Obviously, Conditions 1 and 4 only differ for $n > 1$.

Condition 4 *Given that Condition 3 holds:*

$$\sum_{\substack{l=0 \\ l \neq i}}^n \frac{\partial U_l}{\partial x_i} = 0$$

for $\mathbf{x} = \mathbf{x}^*$ and all $i = 1, \dots, n$.

Proposition 4 *Given that all agents' second order conditions are satisfied, the rotten kid theorem holds for all $W(\mathbf{U})$ if and only if Condition 3 holds and Condition 4 does not hold.*⁹

4 Pareto efficiency

Whereas we have stated the Samaritan's dilemma and the rotten kid theorem in terms of agent 0's first best, an alternative and often used formulation is in terms of Pareto efficiency.⁹ In this section we shall examine the links between the two formulations.

It is clear that agent 0's first best is a Pareto optimum, since any other allocation would make her worse off. The interesting question is: When agent 0 cannot reach her first best when she moves last (first), does this imply that the outcome is not Pareto efficient either?

Let us first establish the relation between agent 0's first best and Pareto efficiency in general. With equation (12) in subsection 3.1, we have already defined the payoff possibility frontier *PPF*. This definition implies:

Lemma 2 *For each allocation \mathbf{U}^* on the Payoff Possibility Frontier (PPF), $dU_k^*/dx_i = 0$ is feasible for all $k = 0, \dots, n$, and all $i = 1, \dots, n$, where dU_0/dx_i and dU_j/dx_i , $j = 1, \dots, n$, are defined by (19) and (20) respectively.*

The idea behind this lemma is the following. Consider a marginal change in x_i , after which agent 0 adjusts \mathbf{y} to compensate all agents j : $dU_j/dx_i = 0$ for all $j = 1, \dots, n$. After this compensation, U_0 should also be back at its original level: $dU_0/dx_i = 0$. Otherwise, U_0 can be increased while all U_j , $j = 1, \dots, n$, remain the same.

⁹Bergstrom ([3], p. 1146) identifies the altruist's first best with "the" Pareto optimum, neglecting the fact that there is a whole range of Pareto optima.

If agent 0 were selfish, then all allocations on the *PPF* would be Pareto efficient. However, when agent 0 is an altruist, a Pareto improvement from some allocations on the *PPF* may be possible. This would be the case if an increase in y_i , which obviously raises U_i , would also increase the value of agent 0's objective function W . Let us now define the Altruistic Payoff Possibility Frontier *APPF* as that part of the *PPF* from which Pareto improvements are impossible:

Definition 5 \mathbf{U}^{A^*} is an element of agent 0's Altruistic Payoff Possibility Frontier (*APPF*) if and only if it is an element of the *PPF* and $dW(\mathbf{U}^{A^*})/dy_i \leq 0$ for all $i = 1, \dots, n$.

Lemma 3 All allocations and only the allocations \mathbf{U}^{A^*} are Pareto efficient.

In our continuous version of the Samaritan's dilemma in subsection 2.1, the sequence where the parasite moves first does not lead to a Pareto optimum. The parasite does not work hard enough, because a higher work effort would decrease the Samaritan's donation. Given the Samaritan's donation, however, the parasite could increase his own payoff and the Samaritan's objective function by working harder. We will now see that this result holds in general. Intuitively, the reason why a Pareto-efficient allocation is not agent 0's first best is that agent 0 does not like the payoff distribution. However, when agent 0 moves last, she determines the payoff distribution. Then, when the allocation is Pareto-efficient, it must be agent 0's first best.

Proposition 5 In the game where agents $i, i = 1, \dots, n$, move in stage one and agent 0 moves in stage two, the outcome is Pareto efficient if and only if it is agent 0's first best.

In our continuous version of the rotten kid theorem (subsection 2.2), when the parent cannot reach her first best when she moves first, the outcome is not Pareto efficient either. Pareto efficiency requires the kid to maximize family income, but the kid will maximize his own income instead. We shall now see that this result can only be generalized partially.

Proposition 6 Consider the game where agent 0 moves in stage one and agents $i, i = 1, \dots, n$, move in stage two. The outcome is Pareto efficient if and only if it is agent 0's first best when:

1. $n = 1$, and/or:
2. agent 0 can reach her first best for all $W(\mathbf{U})$ when agents $i, i = 1, \dots, n$, move in stage one and agent 0 moves in stage two.

Otherwise, the resulting allocation \mathbf{U} is a Pareto optimum if and only if it is on agent 0's APPF according to Definition 5.

5 Bergstrom's rotten kid game

5.1 Introduction

The present paper is not the first to have derived conditions for the rotten kid theorem to hold. Bergstrom [3] and Cornes and Silva [11] have previously derived a condition from a more specific model than ours. In this model, the altruist distributes a certain sum of money among the selfish agents. The total amount of money available may depend on the selfish agents' actions. In this setup, it would be somewhat contrived to study the sequence where the altruist moves first. Accordingly, Bergstrom [3] (at least in his formal model) and Cornes and Silva [11] limit their attention to the positive rotten kid theorem as defined in our Definition 2.

In subsection 5.2, we shall present Bergstrom's [3] solution for the positive rotten kid theorem. We shall see that as his maximization problem for the altruist is a special case of our more general problem, his payoff condition is an accordingly special version of our payoff condition. In subsection 5.3, we shall discuss Cornes and Silva's [11] condition for the positive rotten kid theorem to hold in Bergstrom's [3] model. We shall see that this condition does not carry over to our own more general model and that there are no further solutions to our or Bergstrom's [3] model.

5.2 Bergstrom's solution

In this subsection, we shall discuss Bergstrom's [3] conditions for the positive rotten kid theorem. His condition on the payoff functions is slightly different from our Condition 3. The difference can be traced to differences in the altruist's maximization problem. We shall also discuss the differences in the additional conditions and their derivation.

In Bergstrom’s [3] model, the role of the altruist is limited to the distribution of a certain amount of money. We shall move from our model to Bergstrom’s [3] in two steps. First, let us derive the condition for the positive rotten kid theorem in case agent 0’s objective function includes U_0 , but \mathbf{y} is restricted to “money”. The relevant property of money in this context is the following:

Definition 6 *When \mathbf{y} is money, agent 0’s payoff depends on how much she does on aggregate for all other agents, but not on the distribution of this total amount among the agents. Then the altruist’s payoff $U_0(\mathbf{y}, \mathbf{x})$ is given by $U_0(y_0, \mathbf{x})$ with $y_0 \equiv -\sum_{i=1}^n y_i$.*

Applying Definition 6 of money to our payoff condition (??), we find $\partial U_0/\partial y_0 = \partial U_i/\partial y_i$, which results in payoff functions of the form:

$$U_k = A(\mathbf{x})y_k + B_k(\mathbf{x}) \quad (21)$$

The second and final step from our framework to Bergstrom’s [3] is to exclude U_0 from agent 0’s objective function and to introduce a budget constraint for y_0 . Agent 0’s maximization problem is now:

$$\max W(U_1(y_1, \mathbf{x}), \dots, U_n(y_n, \mathbf{x})) \quad s.t. \quad \sum_{i=1}^n y_i = y(\mathbf{x}) \quad (22)$$

This restriction does not result in a further restriction on the payoff functions. Thus, the positive rotten kid theorem holds for all $W(\mathbf{U})$ and all $y(\mathbf{x})$ if and only if U_k has the form (21) for $k = 1, \dots, n$. This is exactly the condition that Bergstrom [3] derives for the positive rotten kid theorem.

It should be noted that the difference between Bergstrom’s condition (21) and our Condition 3 is irrelevant when agent 0 moves last. This is the sequence that we are primarily interested in. In fact, as we have argued above, this is the only sequence one can meaningfully study in Bergstrom’s maximization problem (22).

The intuition behind the equivalence of conditions (21) and 3 when agent 0 moves last is the following. The positive rotten kid theorem holds if there is only one good, which we might call money. Bergstrom [3] assumes from the outset that \mathbf{y} is money, whereas we have not restricted the nature of \mathbf{y} . However, when the positive rotten kid theorem holds, \mathbf{y} must be money. Thus, the *a priori* restriction of \mathbf{y} to money does not bind.

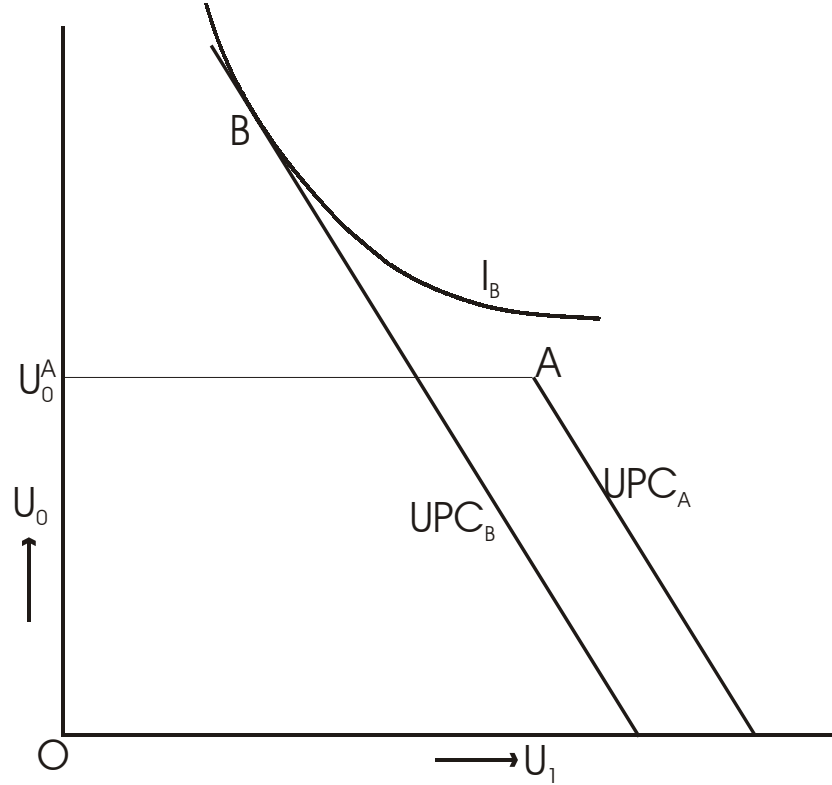


Figure 2: The importance of money

Let us now look at the additional conditions for the positive rotten kid theorem. Bergstrom's [3] Proposition 3 states that when the positive rotten kid theorem holds, all U_i are normal goods and money is important enough, then payoff functions have the form (21). The assumption of normal goods is necessary for the second order conditions to hold in our framework as well, but we have not encountered anything resembling the condition that money is important enough. The rationale behind this condition is illustrated in Figure 2. Consider parallel UPC s with UPC_A the outermost curve. The highest attainable value of U_0 on UPC_A is U_0^A . On UPC_B however, U_0 can rise above U_0^A . Money is not important enough to raise U_0 above U_0^A on UPC_A . In this case, the PPF consists of more than one UPC , although the UPC s are parallel. This creates a problem for the altruist with indifference curve I_B whose first best is point B on UPC_B . Agent 1 will not implement UPC_B , because when he sets UPC_A , agent 0 will choose point A with higher U_1 than point B . We have implicitly excluded this problem by assuming that equilibria are always characterized by internal solutions. In Figure 2, the equilibrium A

of the game is a corner solution, because the altruist’s indifference curve is not tangent to UPC_A .

Finally, there is a difference in approach. Bergstrom [3] mentions the trade-off between restricting W and restricting \mathbf{U} , suggesting that the normal good assumption is a very mild restriction on W . Unlike Bergstrom, we formally derive the payoff condition from the requirement that the first order conditions for the game and the altruist’s first best should coincide.¹⁰

5.3 Cornes and Silva’s solution

In subsection 5.2, we have seen that Bergstrom’s [3] own condition (21) for the positive rotten kid theorem in his game is a special case of our Condition 3. Cornes and Silva [11] recently found another and completely different condition for the positive rotten kid theorem to hold in Bergstrom’s [3] framework. Under this condition, all kids contribute to a pure public good. The reason why Cornes and Silva [11] could find an additional condition is that they, unlike Bergstrom [3], have put a restriction on the budget constraint.

In this subsection we discuss Cornes and Silva’s [11] result in the light of our own analysis, demonstrating why it does not carry over to our more general framework. We shall also argue that there are no additional conditions under which the rotten kid theorem holds for all $W(\mathbf{U})$, neither in Bergstrom’s [3] framework, nor in our more general setup.

In the notation of this paper, Cornes and Silva’s [11] model can be described as follows. Agent i , $i = 1, \dots, n$, only affects the others through his contribution x_i to a pure public good $X \equiv \sum_{i=1}^n x_i$. Agent i has to decide how much x_i of his initial exogenous endowment m_i to contribute to the pure public good. The rest of the endowment plus the transfer t_i from agent 0 is available for consumption y_i of the private good. Agent 0’s budget is zero: $\sum_{i=1}^n t_i = 0$, so that she will also take away from some agents: $t_i < 0$ is feasible. Agent 0’s budget constraint can also be written as $\sum_{i=1}^n y_i = M - X$, with $M \equiv \sum_{i=1}^n m_i$.

The difference between Bergstrom’s [3] and Cornes and Silva’s [11] condition is that Bergstrom’s [3] condition works for all $W(\mathbf{U})$ and all $y(\mathbf{x})$, whereas Cornes and Silva’s

¹⁰Bergstrom [3] restricts the proof of his Proposition 3 to the case of two kids “where the geometry allows an easy, intuitive proof. Extension to higher dimensions is not difficult, but the exposition is tedious.” ([3], p. 1153) We hope our formal exposition is less tedious than the one that Bergstrom [3] had in mind.

[11] condition works for all $W(\mathbf{U})$, but with a restriction $y(\mathbf{x}) = y(\sum_{i=1}^n x_i)$ on the budget constraint. Intuitively, for the rotten kid theorem to hold, agent 0 must react in the same way to any change in \mathbf{x} , setting $dU_k/dx_i = 0$ for all $k = [0, 1], \dots, n$ and all $i = 1, \dots, n$. The Cornes and Silva [11] solution achieves this standardization by defining x_i as agent i 's contribution to the pure public good X . In our solution, the standardization follows from the fact that all agents i contribute to aggregate income I as defined by (??).¹¹

Cornes and Silva [11] only show that the pure public good case is sufficient for the positive rotten kid theorem to hold in Bergstrom's [3] framework. They do not address the issue whether there might still be more solutions. We shall now see that there are no additional solutions to Bergstrom's [3] problem.

First, let us briefly present the derivation of Bergstrom's [3] own solution with our method from subsection 3.3. Analogous to Lemma 1, $dU_j/dx_i = 0$ must hold for all $i, j = 1, \dots, n$ for the rotten kid theorem to apply for all $W(\mathbf{U})$. The agents i set $dU_i/dx_i = 0$ themselves. We need conditions on \mathbf{U} to make sure that agent 0 will set $dU_l/dx_i = 0$ for all other $l, i = 1, \dots, n, l \neq i$. These conditions are (21).

How can we possibly find an additional payoff condition for all $W(\mathbf{U})$? Obviously, this condition should also yield $dU_l/dx_i = 0$ for all $l, i = 1, \dots, n, l \neq i$. In deriving condition (21), we have assumed that agent 0 would have to set all $dU_l/dx_i = 0$ herself. Alternatively, we could impose some restrictions R on the payoff functions $U_i(y_i, \mathbf{x})$ so that $dU_i/dx_i = 0$ automatically implies $dU_l/dx_i = 0$ for some (but not all) $l, i = 1, \dots, n, l \neq i$. However, it can be shown that as long as agent 0 still has to set some $dU_l/dx_i = 0$ herself, the payoff condition will simply be (21) with restrictions R .

The only option left is then to impose that when agent i sets $dU_i/dx_i = 0$, this should automatically imply $dU_i/dx_l = 0$ for all $l, i = 1, \dots, n, l \neq i$. This will be the case if and only if we can define $X \equiv \sum_{i=1}^n x_i$. Then the payoff functions become $U_i(y_i, \mathbf{x}) = U_i(y_i, X)$ and the resource constraint turns into $y(\mathbf{x}) = y(X)$. The n^2 conditions $dU_i/dx_j = 0, i, j = 1, \dots, n$, for implementation of agent 0's first best reduce to n conditions $dU_i/dX = 0$. Agents i 's first order conditions are also $dU_i/dX = 0$.

Without loss of generality, we can specify $y(X) = M - X$. Then we have reproduced

¹¹One could also argue that aggregate income is a public good: All agents $k, k = 0, \dots, n$, benefit from an increase in aggregate income, since all agents i 's, $i = 1, \dots, n$, payoffs are normal goods to agent 0.

Cornes and Silva's [11] pure public good case.

Following the above reasoning, it is clear why Cornes and Silva's [11] condition does not carry over to our more general framework. In the pure public good case where $X \equiv \sum_{i=1}^n x_i$, the agents i , $i = 1, \dots, n$, will set $dU_i/dX = 0$. However, this is not sufficient. We still have to make sure that agent 0 will set $dU_0/dX = 0$. She will do this if and only if the payoff functions satisfy Condition 3 with \mathbf{x} replaced by $X \equiv \sum_{i=1}^n x_i$. Thus, it is impossible to find any solution other than Condition 3 in the general framework.

We conclude that an additional solution for the positive rotten kid theorem can only exist when agent 0 does not have to set any $dU_k/dx_i = 0$ herself. In our general framework this is not feasible, but in Bergstrom's [3] more restricted setup, it is. The additional solution in Bergstrom's [3] setup is exactly Cornes and Silva's [11] pure public good case.

6 Conclusion

For twenty-five years, the Samaritan's dilemma (Buchanan [7]) and the rotten kid theorem (Becker [1] [2]), with their mutually exclusive claims, have coexisted in the economic theory of altruism. This paper has been the first to analyze the conditions on the payoff functions under which either result holds for any altruistic objective function. We have seen that the altruist can reach her first best when she moves first if and only if a selfish agent's action does not affect any other agent's payoff in the optimum. Then there are no externalities to the selfish agents' actions. The altruist can reach her first best when she moves last if and only if there is just one commodity involved. Then the selfish agents cannot manipulate the altruist's trade-off between her own and the selfish agents' payoffs. The selfish agents will maximize the aggregate amount of the single commodity and the altruist will redistribute the commodity.

The theory of altruism can also be applied to government policy. The link between these two fields of research is that the government can be regarded as an altruist, when it maximizes social welfare or any other objective function that depends positively on the payoff of some other player. Thus, the theory of altruism can contribute to our understanding of when collective and individual interests coincide (Shapiro and Petchey [25], Munger [22]). Under the conditions of the Samaritan's dilemma, the government can

reach the optimum if and only if it can commit to a certain policy. If the Samaritan's dilemma does not apply, commitment does not result in the first best. The government may then be better off with a time-consistent policy. Under the conditions of the rotten kid theorem, the time-consistent policy even results in the first best. Starting with Kydland and Prescott [17], most analyses of time consistency have used a more complicated setup than ours.¹² However, a general framework for the study of time consistency issues is still lacking. Our simple model of altruism would be a useful starting point for the development of such a framework (Dijkstra [14]).

7 Appendix

Proof of Lemma 1. Since agent 0 moves last, the first order conditions (10) for agent 0's first best with respect to \mathbf{y} are satisfied. Substituting (10) into the first best conditions (11) for \mathbf{x} , we can rewrite them as:

$$\sum_{k=0}^n W_k \frac{dU_k}{dx_i} = 0 \quad (23)$$

for all $i = 1, \dots, n$, where dU_k/dx_i , $k = 0, \dots, n$, is defined by (19) and (20).

In the equilibrium of the game, agent i , $i = 1, \dots, n$, sets $dU_i/dx_i = 0$. This will result in the first best condition (23) for all $W_k > 0$, $k = 0, \dots, n$, if and only if $dU_i/dx_i = 0$ implies $dU_l/dx_i = 0$ for all $i = 1, \dots, n$, $l = 0, \dots, n, l \neq i$ in agent 0's first best, characterized by $\mathbf{x} \in \mathbf{X}^*$. This is Condition 2.

Proof of Proposition 1. Combining Condition 1 with agents i 's first order conditions for the maximization of U_i (14), we obtain the first best conditions for \mathbf{x} (11). Substituting Condition 1 into agent 0's first order conditions for the maximization of W (15), we obtain the first best conditions for \mathbf{y} (10). This proves the "if" part. The "only if" part follows from the requirement that (14) and (15) should turn into (11) and (10) for all values of $W_k > 0$. This is only possible when Condition 1 holds.

Proof of Proposition 2. By Lemma 1, $dU_k/dx_i = 0$, $k = 0, \dots, n$, $i = 1, \dots, n$, must hold in agent 0's first best. Substituting this into the derivative of agent 0's reaction

¹²Dijkstra [13] offers a straightforward application of the Samaritan's dilemma to time consistency.

function (18), the first and third term on the LHS of (18) drop out. This leaves:

$$W_0 \frac{d(\partial U_0 / \partial y_j)}{dx_i} + W_j \frac{d(\partial U_j / \partial y_j)}{dx_i} = 0$$

Substituting agent 0's first order conditions (16) for \mathbf{y} , this becomes:

$$\frac{d\left(\frac{\partial U_0}{\partial y_j}\right) / dx_i}{\partial U_0 / \partial y_j} = \frac{d\left(\frac{\partial U_j}{\partial y_j}\right) / dx_i}{\partial U_j / \partial y_j} \quad (24)$$

Agent 0 can reach her first best for all $W(\mathbf{U})$ when she moves after agents i if and only if (24) holds for all $i, j = 1, \dots, n$ in agent 0's first best. For a given $x \in X^*$, define the Utility Possibility Contour as the set of all feasible \mathbf{U} . The slope of the *UPC* in dimension j represents the tradeoff between U_0 and U_j :

$$\frac{dU_0}{dU_j} \equiv \frac{\partial U_0 / \partial y_j}{\partial U_j / \partial y_j}$$

Condition (24) implies that a marginal change in x_i , which leads to a different *UPC*, produces a *UPC* with the same slope:

$$\frac{d(dU_0 / dU_j)}{dx_i} = 0$$

Since all *UPCs* on the *PPF* are parallel, *UPCs* on the *PPF* do not intersect, and there can only be one *UPC* implementing the whole *PPF*.

Proof of Proposition 5. The “if” part is obvious. With respect to the “only if” part, note that agent 0 maximizes W with respect to \mathbf{y} according to (10) in stage two:

$$W_0 \frac{\partial U_0}{\partial y_j} + W_j \frac{\partial U_j}{\partial y_j} = 0$$

By Lemmas 2 and 3, there exist dy_j/dx_i for all $i = 1, \dots, n$ such that a Pareto optimum satisfies:

$$W_0 \left(\frac{\partial U_0}{\partial x_i} + \sum_{j=1}^n \frac{\partial U_0}{\partial y_j} \frac{dy_j}{dx_i} \right) + \sum_{j=1}^n W_j \left(\frac{\partial U_j}{\partial x_i} + \frac{\partial U_j}{\partial y_j} \frac{dy_j}{dx_i} \right) = 0$$

Substituting (10), this becomes:

$$\sum_{k=0}^n W_k \frac{\partial U_k}{\partial x_i} = 0$$

These are (11), the first order conditions for W with respect to \mathbf{x} . Thus, all first best conditions are satisfied and the allocation is agent 0's first best.

Proof of Proposition 6. In Case 1, by Lemmas 2 and 3, \mathbf{U} can only be Pareto efficient if $dU_0/dx = dU_1/dx = 0$ is feasible. For dU_1/dx , we find:

$$\frac{dU_1}{dx} = \frac{\partial U_1}{\partial x} + \frac{\partial U_1}{\partial y} \frac{dy}{dx} = \frac{\partial U_1}{\partial y} \frac{dy}{dx}$$

The second equality follows from the fact that agent 1 sets $\partial U_1/\partial x = 0$ in stage two. Thus, $dU_1/dx = 0$ implies $dy/dx = 0$. Then for dU_0/dx :

$$\frac{dU_0}{dx} = \frac{\partial U_0}{\partial x} + \frac{\partial U_0}{\partial y} \frac{dy}{dx} = \frac{\partial U_0}{\partial x}$$

Thus, $dU_0/dx = dU_1/dx = 0$ is feasible if and only if $\partial U_0/\partial x = 0$. But then Condition 1 is satisfied and \mathbf{U} is agent 0's first best.

In Case 2, Condition 3 holds by Proposition 2. This means that the \mathbf{x}^* that implements agent 0's first best implements the whole PPF . When agent 0 moves first, $\mathbf{x} \neq \mathbf{x}^*$, because she cannot reach her first best. Then the allocation is not on the PPF . By Lemma 3, when an allocation is not on the PPF , it is not Pareto efficient either.

References

- [1] Becker, Gary S. (1974), "A theory of social interaction", *Journal of Political Economy* 82: 1063-1093.
- [2] Becker, Gary S. (1976), "Altruism, egoism, and genetic fitness: Economics and sociobiology", *Journal of Economic Literature* 14: 817-826.
- [3] Bergstrom, Theodore C. (1989), "A fresh look at the rotten kid theorem— and other household mysteries", *Journal of Political Economy* 97: 1138-1159.
- [4] Bernheim, B. Douglas, Andrei Schleifer and Lawrence H. Summers (1985), "The strategic bequest motive", *Journal of Political Economy* 93: 1045-1076.
- [5] Bruce, Neil and Michael Waldman (1990), "The rotten-kid theorem meets the Samaritan's dilemma", *Quarterly Journal of Economics* 105: 155-165.

- [6] Bruce, Neil and Michael Waldman (1991), "Transfers in kind: Why they can be efficient and nonpaternalistic", *American Economic Review* 81: 1345-1351.
- [7] Buchanan, James M. (1975), "The Samaritan's dilemma", in: E.S. Phelps (ed.), *Altruism, Morality, and Economic Theory*, Sage Foundation, New York, 71-85.
- [8] Coate, Stephen (1995), "Altruism, the Samaritan's dilemma, and government transfer policy", *American Economic Review* 85: 46-57.
- [9] Chami, Ralph (1996), "King Lear's dilemma: Precommitment versus the last word", *Economics Letters* 52: 171-176.
- [10] Chami, Ralph (1998), "Private income transfers and market incentives", *Economica* 65: 557-580.
- [11] Cornes, Richard C. and Emilson C.D. Silva (1999), "Rotten kids, purity, and perfection", *Journal of Political Economy* 107: 1034-1040.
- [12] Cox, Donald (1987), "Motives for private income transfers", *Journal of Political Economy* 95: 508-546.
- [13] Dijkstra, Bouwe R. (2000), "Investment incentives of environmental policy instruments", Discussion Paper 308, Faculty of Economics, University of Heidelberg.
- [14] Dijkstra, Bouwe R. (2001), "Policy commitment vs. time consistency", Interdisciplinary Institute for Environmental Economics, University of Heidelberg.
- [15] Hirshleifer, Jack (1977), "Shakespeare vs. Becker on altruism: The importance of having the last word", *Journal of Economic Literature* 15: 500-502.
- [16] Jürges, Hendrik (2000), "Of rotten kids and Rawlsian parents: The optimal timing of intergenerational transfers", *Journal of Population Economics* 13: 147-157.
- [17] Kydland, Finn E. and Edward C. Prescott (1977), "Rules rather than discretion: The inconsistency of optimal plans", *Journal of Political Economy* 85: 473-491.

- [18] Lagerlöf, Johan (1999), “Incomplete information in the Samaritan’s dilemma: The dilemma (almost) vanishes”, Discussion Paper FS IV 99-12, Wissenschaftszentrum Berlin.
- [19] Lindbeck, Assar and Jörgen W. Weibull (1988), “Altruism and time consistency: The economics of fait accompli”, *Journal of Political Economy* 96: 1165-1182.
- [20] Lord, William and Peter Raganzas (1995), “Uncertainty, altruism, and savings: Precautionary savings meets the Samaritan’s dilemma”, *Public Finance* 50: 404-419.
- [21] Monderer, Dov and Lloyd S. Shapley (1996), “Potential games”, *Games and Economic Behavior* 14: 124-143.
- [22] Munger, Michael C. (2000), “Five questions: An integrated research agenda for Public Choice”, *Public Choice* 103: 1-12.
- [23] Pollak, Robert A. (1985), “A transactions cost approach to families and households”, *Journal of Economic Literature* 23: 581-608.
- [24] Schmidtchen, Dieter (1999), “To help or not to help: The Samaritan’s dilemma revisited”, Discussion paper 9909, Center for the Study of Law and Economics, Saarland University.
- [25] Shapiro, Perry and Jeffrey Petchey (1998), “The coincidence of collective and individual interests”, Working paper 9-98R, Department of Economics, University of California, Santa Barbara.
- [26] Wigger, Berthold U. (1996), “Two-sided altruism, the Samaritan’s dilemma, and universal compulsory insurance”, *Public Finance* 51: 275-290.