# TIME CONSISTENCY AND RELATED CONCEPTS

Daniel B. Klein
Department of Economics
University of California
Irvine, California 92717

Brendan O'Flaherty
Department of Economics
Columbia University
New York, New York 10027

Revised: March 1993

# Time Inconsistency and Related Concepts

by

Daniel Klein and Brendan O'Flaherty

## Abstract

We give rigorous game-theoretic meaning to the Stackelberg notions of time inconsistency and to the idea of commitment being of value (or "sequential irrationality"). Time inconsistency treats desirable deviations only along the path, whereas sequential irrationality treats deviations everywhere in the tree.

We establish four relationships: time inconsistency implies sequential irrationality; sequential irrationality implies and is implied by the failure of subgame perfect equilibrium; and time consistency implies Nash equilibrium. (In the last case the restrictions are numerous.) We show also that no other relationships are valid. Our basic results are easily gleaned from the schema provided in Figure 1 (which follows typescript page 28).

# 1. INTRODUCTION

Ever since Kydland & Prescott [1977] showed that rules could be better than discretion, economists interested in how governments behave have turned to game theory to try to understand the threats and promises that governments might make. The subject is commonly referred to as time (or dynamic) inconsistency, and has become an important part of macroeconomics, international trade theory, and public finance. This grafting of noncooperative game theory on to the more policy-oriented fields has generally been quite productive.

The grafting, however, has not been perfect. There are subtle differences that have often been overlooked between noncooperative game theory and the policy-oriented traditions. The purpose of this paper is to show what these differences are. Our investigation is based on two distinctions: first, the distinction between the Nash solution concept and the Stackelberg solution concept; second, the distinction between deviations "along the path" and deviations "anywhere in the tree."

## The Distinction Between Nash and Stackelberg

With a Nash solution concept, everyone is a follower. Nash equilibrium and its refinements require best strategies given the strategy choices of the other players. Writers from the noncooperative game theory tradition are used to the Nash solution concept.

The Stackelberg solution concept recognizes one player as the leader (or, in our terminology, "ruler"), who pre-emptively chooses a strategy and thereby influences the strategy choices of the other players. In Stackelberg equilibrium, the followers are doing as well as possible given the ruler's strategy, while the ruler is doing as well as possible given the followers' response function. Of course any Stackelberg game can be translated into a Nash counterpart, but much may be lost in the translation.[1] Writers from public finance, international trade, and macroeconomics find the Stackelberg framework natural, with the government acting as ruler.

## The Distinction Between Deviations Along-the-Path and Anywhere-in-the-Tree

In Nash games, Nash equilibrium requires best choices only along the path, whereas subgame perfection requires best choices in all subgames. We make a similar distinction for Stackelberg games. If, along the path of an optimal plan, the ruler has no desire to deviate to a new plan, then the original plan is time consistent. If the ruler would have no such desire *anywhere* in the tree, we say the plan is sequentially rational. "Sequential rationality" -- the Stackelberg analog to subgame perfection -- is a concept that has not been formally identified in the literature previously. (Our usage of the term is distinct from what Kreps & Wilson [1982] mean by it.) This notion appears to be very useful. It is probably implicit in many discussions of public finance and macroeconomics over the last decade. Sequential irrationality is nearly synonymous with commitment being of value to the ruler.

Notice that, unlike in the Nash test, when a Stackelberg ruler deviates to a new plan, subsequent follower choices are revised as well. That is, the leadership quality of the ruler holds for deviations as well as for the original

announcement.

## Depiction of Our Results

Using the two distinctions we create the following matrix:

[Please consult Figure 1.]

The arrow pointing from "subgame perfect equilibrium" to "Nash equilibrium" means that the former implies the latter, which is well known (Selten [1975]). In this paper we establish the other four arrows shown. The dashed arrow from "time consistency" to "Nash equilibrium" signifies that nontrivial restrictions are required for the relationship. (Each of our results requires a condition on ties.) The paper establishes also that no arrows can be added to the matrix.

Fershtman [1989] and others have discussed the relationship between subgame perfect equilibrium and time consistency, but only in special cases.

We do not assume that the ruler is benevolent, as is typically done in the time inconsistency literature. Ruler benevolence is a special case in our investigation.

The next section sets out notation. The notation is elaborate because we need to fit Stackelberg notions into a noncooperative game theory fabric. This has not been done previously. In doing so, various fine points arise that require fresh attention. Section 3 presents our results. There is a theorem for each of the arrows shown in Figure 1. Once precisely stated, the theorems corresponding to the solid arrows are easily proved. The theorem corresponding to the dashed arrow requires a lengthy proof, which is in

Appendix 1.

## 2. NOTATION AND DEFINITIONS

This section provides most, but not all, of the definitions used in the paper. The exposition of this section is organized in numbered subsections.

(2.1) An <u>S-game</u> is a four-tuple $\Sigma = (G, i, C, s)$.

(2.2) <u>The reference game</u> G is an extensive form game.

(2.3) <u>The ruler, i, and the public</u>. Let I denote the index set of players of G, where $|I| = m+1$. Player i is called the <u>ruler</u>, and the set $I \setminus \{i\}$ of other players is called the <u>public</u>. For the remainder of this paper we let player 1 be the ruler (or i=1). The set of nodes assigned to the ruler is R; the set of nodes assigned to the public is P.

(2.4) <u>Information assumption on G</u>. We make the following three assumptions about G: (a) every ruler information set is a singleton, and (b) at every information set in G all previous play by the ruler has been publicly observed. Assumptions (a) and (b) are rather natural since the ruler announces her plan in advance of play and she knows how the public will respond. The assumptions ensure that a subgame originates at every ruler node.

(2.5) <u>Plans</u>. Let B' denote the set of player i's behavior strategies in G. In

the S-game $\Sigma$, a subset B of B' is the set of permissable <u>plans</u> for the ruler.
(Subsection (2.14) suggests why the ruler may not have access to the complete
set of behavior strategies, B'.) Think of a plan as an announcement of what
the ruler will do at each of her nodes.

(2.6) <u>Public response function, s</u>. Let S denote the set of behavior strategy
m-tuples that exist in G for the public. Let s(b), the <u>public response function</u>,
denote the unique element of S that is picked out when plan b is announced.
For the remainder of the paper we assume that for every plan b s(b) is a
subgame perfect equilibrium to the m-player game induced by b.

(2.7) <u>Ruler payoffs</u>. U(b, d) denotes the ruler's payoff when she uses plan b
and the public uses behavior strategy m-tuple $d \in S$. For U(b, s(b)) we will
sometimes employ the summary payoff function u(b) := U(b, s(b)).

(2.8) <u>Optimal plans</u>. A plan b* is <u>optimal</u> iff for all $b \in B$, $u(b^*) \geq u(b)$.

(2.9) <u>Local variations and subgames</u>. For any node $v \in R$, let $b_v$ denote the
local strategy specified by b at v. Let $b \mid b'_v$ denote the ruler's behavior strategy
that results if the local strategy assigned by b to node $v \in R$ is changed to $b'_v$
while the local strategies assigned by b to other nodes in R remain
unchanged.

For any node t where a subgame originates let G(t) denote the subgame

whose origin is t. Let b(t) denote the behavior strategy induced by plan b on

G(t), and denote the ruler's payoff function on G(t) as $U_t(.)$, and her summary

payoff function as $u_t(.)$. Let s(b(t)) denote the behavior strategy m-tuple

induced on G(t) by s(b).

(2.10) <u>Sequential rationality</u>. A plan b is <u>sequentially rational</u> iff for every

v ε R and every local strategy $b'_v$ at v

$$u_v( b(v) ) \geq u_v( b(v) | b'_v).  \qquad (1)$$

To understand this definition think about applying (1) backward through the

game. If a plan is not sequentially rational, it is <u>sequentially irrational</u>.

Consult Figures 2 ("If you scratch my back I'll scratch your back") and

Figure 3 ("Scratch my back or else I'll break your back") for simple examples

of sequentially irrational plans.

[Please consult Figures 2 & 3.]

Using the idea of sequential irrationality and plan optimality, we get a

straightforward and natural standard for whether the ruler values

commitment conveyance. Without commitment conveyance the ruler is

restricted to sequentially rational plans. Therefore, were the ruler to have

commitment conveyance and were all her optimal plans to be sequentially

irrational, then we say she benefits from having commitment conveyance.

We call this property "commitment dominance." Thus, when all optimal

plans are sequentially irrational, the ruler faces commitment dominance.

The commitment interpretation is apparent in Figures 2 and 3. In those

Figures, were the ruler to lack commitment conveyance, her promises and

threats would not be credited, and her back would not be scratched.

(2.11) <u>Time consistency</u>. The <u>concatenated behavior strategy for the ruler</u>

κ(b, v, b') is the behavior strategy that results from behavior strategy b if the

behavior strategy induced by b on G(v), v ε R, is changed to b'(v) while the local

strategies assigned by b to other nodes in R remain unchanged. Similarly, the

<u>concatenated behavior strategy m-tuple for the public</u> γ(d, v, d') is the behavior

strategy m-tuple that results from behavior strategy m-tuple d if the behavior

strategy m-tuple induced by d on G(v), v ε R, is changed to d'(v) while the local

strategies assigned by d to other vertices in P remain unchanged.

A plan b is <u>time inconsistent</u> iff there exists a v ε R and some sequentially

rational plan b' such that

$$U(\ \kappa(b, v, b'),\ \gamma(\ s(b), v, s(b')\ )\ )\ )\ >\ u(b). \tag{2}$$

A plan that is not time inconsistent is called <u>time consistent</u>.[2] Consult Figure

2 for a simple example of a time inconsistent plan.

Our definition of time inconsistency is the natural and faithful

game-theoretic representation of what that term has always meant. The only

distinctive feature of our definition is that, were the ruler able to fool the public

and deviate, at a reached node the ruler can reannounce *only a sequentially

rational plan*. Alternatively one may wish to permit her to reannounce

convincingly any plan, and to fool the public repeatedly. The issue of the

proper choice set at the point of deviation has scarcely arisen in the time

inconsistency literature because in those models the government typically reaches a point of time inconsistency at final moves, so sequential rationality would be its preferred deviation even if it had a wider choice. Our decision to restrict reannouncements to sequentially rational plans conforms to the proverb: "Fool me once, shame on you; fool me twice, shame on me." Once fooled, the public will not believe anything but a sequentially rational plan, which will truly stick for the remainder of the game. We are not wedded to this formulation; other conventions on this issue are worthy of exploration.

(2.12) <u>Historical independence</u>. This is a condition treating ties in citizen payoffs, needed for results concerning Nash and subgame perfect equilibrium. A public response function s satisfies <u>historical independence</u> iff whenever b and b' are two plans that differ on no node that succeeds v ε R, then at every node w ε P that succeeds v the local strategy induced by s(b) is the same as the local strategy induced by s(b'). Expressed again, historical independence is satisfied iff the following is true for any node v ε R and any pair of plans b and b': If $b(v) = b'(v) | b_v$, then $s(b(v)) = s(b'(v))$.

Intuitively, suppose we were to intrude on an S-game at some citizen node w, and that plans b and b' are the same in G(w) (that is, b(w)=b'(w)). Suppose we got on a soapbox at w and said to the citizens active in G(w), "Plans b and b' both imply the same thing from here on. Is anyone of you going to base your response to this on whether b or b' is being used elsewhere?" Historical independence requires a universal answer of no. Because of the local best

reply restriction already placed on s, historical independence can matter only when some citizen faces a tie.

Figure 4 provides an example of a violation of historical independence. Notice that the public can "threaten" the ruler by choosing an appropriate response rule. Considering the second response rule described in Figure 4, the public response of (B) to (R) violates the spirit, if not the letter, of the Stackelberg tradition. (Nalebuff and Shubik [1988] demonstrate the point that ties give players the kind of freedom to make threats and promises that we intuitively associate with the ability to make commitments.)

[Please consult Figure 4.]

(2.13) SR-equivalence. This is a condition treating ties in ruler payoffs. An S-game $\Sigma$ satisfies SR-equivalence iff whenever b and b' are sequentially rational plans

$$u_V(b(v)) = u_V(b'(v))$$

for all $v \in R$.  Not all S-games satisfy SR-equivalence; the game in Figure 5 provides an example.[3]  SR-equivalence can fail only when the ruler faces a tie in payoffs, but ties need not cause failure of SR-equivalence.

[Please consult Figure 5.]

Notice that if the ruler moves only once (that is, every ruler node is without predecessor or successor ruler nodes), SR-equivalence is trivially satisfied for any public response function s that satisfies the local best reply restriction. S-games with this characteristic are the games most heavily studied in the time inconsistency literature on monetary and public economics (e.g.,

Kydland & Prescott [1977]; Fischer [1980]).

(2.14) <u>Policy formation, C</u>. In the Kydland & Prescott (1977) Phillips curve model, a government commitment to inflation policy cannot take the form of any reaction function (with domain being the history of citizen employment decisions). In the time inconsistency literature, the government is restricted to plans that take the form of a single magnitude that applies uniformly across all possible citizen histories standing at the government's "time to act." Hence the literature makes frequent use of the term "open loop" in describing the "rules," or commitment, regime.

Item C -- "policy formation" -- addresses such restrictions on the set of plans the ruler can choose from. We shall not exposit policy formation in this paper because it does not play a role except in Theorem 4. Let us say only that if there are no policy formation restrictions on the ruler's set of eligible plans (that is, B = B'), we say that the S-game has "perfect policy formation." Perfect policy formation is necessary in Theorem 4. A fuller discussion of policy formation is provided in Klein & O'Flaherty (1992), where it plays a central role; for fuller interpretation, see Klein & O'Flaherty (forthcoming).

## 3. RESULTS ABOUT PLANS

Although it is optimal plans that rulers care about, it is desirable to state results as generally as possible. Our results are not confined to optimal plans.

Theorem 1 confirms our sensibility that time inconsistency precludes

sequential rationality, or, in other words, that time inconsistency depends on commitment conveyance.

**Theorem 1:** *Let $\Sigma$ satisfy SR-equivalence.  If a plan b is time inconsistent, then it is sequentially irrational.*

**Proof:** Since b is time inconsistent, there is some $v \varepsilon R$ and some sequentially rational plan b' such that

$$U(\ \kappa(b,v,b'),\ \gamma(\ s(b),v,s(b')\ )\ )\ >\ U(b,\ s(b))$$

Except in G(v) the play of the LHS is identical to the play of the RHS.  The inequality results from differences in G(v):

$$U_v(\ b'(v),\ s(\ b'(v)\ )\ )\ >\ U_v(\ b(v),\ s(\ b(v)\ )\ )$$

Since b' is sequentially rational, by SR-equivalence b must not be sequentially rational.  Hence b is sequentially irrational.

<div align="right">♦ ♦ ♦</div>

SR-equivalence is a necessary restriction on the claim that every inconsistent plan is sequentially irrational, as shown by plan b in Figure 5.

How about the converse of Theorem 1?  That is, does time consistency imply sequential rationality?  The answer is no.  Figure 3 shows an optimal

plan that is time consistent and sequentially irrational.

Theorems 2 through 4 relate plan properties to properties of equilibria in the reference game, linking Stackelberg concepts to familiar concepts in noncooperative game theory. More specifically, Theorems 2 and 3 confirm our intuition about sequential irrationality (or commitment dominance) and the failure of subgame perfect equilibrium.

**Theorem 2:** *Let $\Sigma$ satisfy historical independence and let b be some plan. If b is sequentially irrational, then (b, s(b)) is not a subgame perfect equilibrium of G.*

**Proof:** Suppose b is sequentially irrational. Then there is some $v \in R$ and some $b'_V$ such that

$$U_V( b(v) \mid b'_V, s(b(v) \mid b'_V) ) > U_V( b(v), s(b(v)) ). \tag{3}$$

By historical independence, $s(b(v) \mid b'_V) = s(b(v))$. Substituting into the LHS of (3) yields:

$$U_V( b(v) \mid b'_V, s(b(v)) ) > U_V( b(v), s(b(v)) ) \tag{4}$$

Hence b(v) is not a best response to s(b(v)) in G(v); hence (b, s(b)) is not a subgame perfect equilibrium of G.

$$\blacklozenge\blacklozenge\blacklozenge$$

**Theorem 3:** *Let $\Sigma$ satisfy historical independence and let b be some plan. If (b, s(b)) is not a subgame perfect equilibrium of G, then b is sequentially irrational.*

**Proof:** We have assumed that s(b) is a subgame perfect equilibrium in the citizen game induced by b. Thus if (b, s(b)) is not a subgame perfect equilibrium of G there must be some $v \in R$ and some $b'_v$ such that

$$U_v(\; b(v) \mid b'_v,\; s(b(v))\; ) \; > \; U_v(\; b(v),\; s(b(v))\; ). \tag{5}$$

By historical independence , $s(b(v)) = s(b(v) \mid b'_v)$. Substituting into the LHS of (5) we get,

$$U_v(\; b(v) \mid b'_v,\; s(b(v) \mid b'_v)\; ) \; \geq \; U_v(\; b(v),\; s(b(v))\; ).$$

Hence b is sequentially irrational.

$$\blacklozenge\blacklozenge\blacklozenge$$

Is historical independence necessary for Theorems 2 and 3? The game in Figure 4 and its accompanying discussion show that it is.

From Theorems 1 and 2 we have,

**Corollary 1:** *Let $\Sigma$ satisfy historical independence and SR-equivalence. If a plan b is time inconsistent, then (b, s(b)) is not a subgame perfect equilibrium of G.*

Fershtman [1989] obtained a similar result, but in the context of Markov games, rather than in the context of S-games.

Even with both historical independence and SR-equivalence, neither the converse of Theorem 1 (nor the converse of Corollary 1) holds. As Fershtman argues, time consistency requires that the ruler's plan choose local best replies only along the path of play, while perfection and sequential rationality require that local best replies be chosen at every node. This point is also made by McTaggart & Salant [1990] and by Guiso and Terlizzese [1990]. Not surprisingly, Schelling [1960, 177] first exposited the basic insight: "[A] promise [think time inconsistency] is different from a threat [think subgame imperfection]. The difference is that a promise is costly when it succeeds, and a threat is costly when it fails. A successful threat is one that is not carried out [whereas a successful (and genuine) promise is carried out]."

Now, how about time consistency and Nash equilibrium? Figure 2 shows a time inconsistent plan b with (b, s(b)) not a Nash equilibrium. Figure 3 shows a time consistent plan b with (b, s(b)) a Nash (but not subgame perfect) equilibrium. These figures suggest a relationship between time consistency and Nash equilibrium. Theorem 4 will establish a relationship, but first more refinement is necessary.

Consider Figure 6. Optimal plan b, shown by the arrows, is time consistent, but (b, s(b)) is not a Nash equilibrium. Given b and s(b), the move at k is strategically irrelevant, since no matter what is specified at k, k will not be reached.

[Figure 6 here.]

Formally, for a given S-game with perfect policy formation and a plan b, the move at some $v \in R$ is __strategically irrelevant under b__ iff for any local strategy $b'_v$ at v:

$$\rho(v, (b \mid b'_v, s(b \mid b'_v))) = 0,$$

where $\rho(v, (\beta))$ denotes the realization probability of node v when the players use the behavioral strategy (m+1)-tuple $\beta$. Now, if a move is strategically irrelevant under some plan b, a reasonable condition on b is that the local behavior strategy at that move be locally best. Formally, a plan b is __nonidiosyncratic__ iff for every $v \in R$ such that v is strategically irrelevant under b, and for every $b'_v$

$$U_v( b(v), s(b(v)) ) \geq U_v( b(v) \mid b'_v, s(b(v) \mid b'_v) ).$$

Notice that in Figure 6 b fails this condition at node k and therefore is idiosyncratic. Nonidiosyncrasy is one condition we will place on b in Theorem 4.

Figure 7 shows that we need yet another condition for our relationship between time consistency and Nash equilibrium. Plan b is shown by the arrows and specifies 0.5 probability on each choice at x. For time inconsistency we permit the ruler to revert to sequentially rational plans only, and at x the ruler would not want to revert to the sequentially rational plan. Thus b is time consistent, as well as nonidiosyncratic, yet (b, s(b)) is not a Nash equilibrium. This result is avoided if we require also that b be a pure strategy, which is our final condition for Theorem 4.

[Figure 7 here.]

**Theorem 4:** *Let $\Sigma$ satisfy historical independence (and perfect policy formation), and let b be some plan. If b is time consistent, nonidiosyncratic, and pure, then (b,s(b)) is a Nash equilibrium of G.*

**Proof:** The proof is in Appendix 1.

Theorem 4 says that with several restrictions time consistency implies Nash equilibrium. The proof shows primarily that b must be a best response to s(b). Consider the following: Suppose b were not a best response to s(b) whereas b' were. The differences in local strategies between b and b' cannot be limited to a single vertex because then b would be time inconsistent. But necessary changes in local strategies of b' cannot be more than one since, while the first change opens up a new direction, it leads to nodes that are strategically irrelevant under b, so the b moves down that road are already best. Thus no superior b' can exist, and b must be a best response to s(b). (A

simple example showing the necessity of the perfect policy formation condition is available from the authors.)

Is the converse true? That is, assuming whatever minor restrictions may be necessary, if (b, s(b)) is a Nash equilibrium, must b be time consistent? The answer is no, as shown by Figure 8.

[Figure 8 here.]

## 4. CONCLUDING COMMENTS

Borrowing machinery from the extensive form games, we have given time consistency rigorous game-theoretic expression. We establish a relationship between time consistency and Nash equilibrium, a relationship between time inconsistency and sequential irrationality, and converse relationships between sequential rationality and subgame perfect equilibrium. Our basic results are easily gleaned from the schema in Figure 1.

The time inconsistency literature underscores two points: first, that commitments can help the ruler; second, that the ruler may get to point at which she will want to change plans. Sometimes in print and often in casual discussions these two ideas have been used interchangeably. Perhaps the frequency of this error can be explained in the following manner: An "open-loop" programing approach to economic problems is essentially a Stackelberg approach. When an optimal plan is time consistent there is no way of knowing whether the plan outperforms sequentially rational plans (or, equivalently, whether the ability to convey a commitment is valuable) *except by making an explicit comparison*. When the optimal plan is time inconsistent,

however, in all but a rather unimportant fringe of cases, the investigator

knows that the plan violates sequential rationality and therefore that

commitment is valuable. That is why need for commitment conveyance has

been closely -- sometimes too closely -- associated with time inconsistency.

**Appendix 1:** Proof of Theorem 4.

To prove Theorem 4 we will use Lemma 1.

**Lemma 1:** *Let b be some plan. Let d be any perfect best response to s(b), that is,*

*for every v $\varepsilon$ R and for every local strategy d'$_v$ at v,*

$$U_v(\ d(v),\ s(b(v))\ )\ \geq\ U_v(\ d(v)\,|\,d'_v,\ s(b(v))\ ).$$

*Let g be the perfect best response to s(b) that results from replacing the local strategy of d with local strategy of b at each v where*

$$U_v(\ d(v)\,|\,b_v,\ s(b(v))\ )\ =\ U_v(\ d(v),\ s(b(v))\ ). \tag{6}$$

$$Let\ Y := \left\{v\ \varepsilon R:\ U_v(\ g(v),\ s(b(v))\ )\ >\ U_v(\ g(v)\,|\,b_v,\ s(b(v))\ )\right\}. \tag{7}$$

$$Let\ X := \left\{v\ \varepsilon R:\ \rho(v,\ (g,\ s(b)))\ >\ 0\right\}. \tag{8}$$

*If X $\cap$ Y is empty, then b is a best response to s(b).*

**Proof:** Suppose X$\cap$Y is empty. If X is empty, citizens begin the game and s(b) terminates the play before the ruler even makes a move. Clearly any strategy by the ruler is a best response to s(b).

For X nonempty, consider each v $\varepsilon$ X such that v has no predecessors in R. Each such v is not in Y by supposition, so (6) must be satisfied there and the

local strategy from b is specified by g. Consider the "next level" of v $\varepsilon$ X (that is, successor ruler vertices in X such that no other ruler vertices come between). Again, Y is not satisfied there so the local strategy from b is

specified by g at each such ruler vertex. Continuing through the game by considering successive levels of ruler vertices in X, we see that the path of (g, s(b)) is identical to that of (b, s(b)). Therefore

$$U(g, s(b)) = U(b, s(b)).$$

Since g is a best response to s(b), b must also be a best response to s(b),

♦♦♦

**Theorem 4:** *Let $\Sigma$ satisfy historical independence and perfect policy formation, and let b be some plan. If b is time consistent, nonidiosyncratic, and pure, then (b, s(b)) is a Nash equilibrium of G.*

**Proof:** By the supposition that S(b) assigns local best replies, each citizen is playing a best response to the stategies of all the other players. We need to show that b is a best response to s(b).

Let g, Y, and X be as defined in Lemma 1. We proceed to show that X∩Y is empty.

Let W be the origin of G (that is, W has no predecessors in G). Letting u represent any vertex/terminal node in G, let C(W, u) be the <u>curve</u> connecting W and u. A strategy is said to specify a <u>deviation</u> from C(W, u) if it places zero probability on at least one move along C(W, u).

Divide R (the set of ruler vertices) into three exhaustive disjoint sets:

$R_1 := \{v \in R: \rho(v, (b, s(b))) > 0\}$

$R_2 := \{v \in R: b \text{ specifies at least one deviation from } C(W, v)\}$

$R_3 := \{v \in R: b \text{ does not specify a deviation from } C(W, v) \text{ but } s(b) \text{ does}\}$

Where u and n are two vertices in G, the notation u < n means that u

precedes n. From Y we can generate sets $Y_i$ (i=1,2,...,z) of ordered vertices:

$$Y_i := \{\, v_1, v_2,..., v_{\omega i} \in Y: v_1 < v_2 < ... < v_{\omega i} \,\}$$

The $Y_i$'s certainly may have a nonempty intersection, although a common ruler vertex may have different tags depending on the i.

<u>Claim(i): No $v_\omega$ is in $R_1$.</u> For any $Y_i$, consider the last vertex in the set, $v_{\omega i}$. For convenience we will drop the subscript i on on $v_{\omega i}$. Since no element of Y succeeds $v_\omega$,

$$U_{v\omega}(\, g(v_\omega), s(b(v_\omega))\,) = U_{v\omega}(\, b(v_\omega) \mid g_{v\omega}, s(b(v_\omega))\,), \tag{9}$$

and that

$$U_{v\omega}(\, b(v_\omega), s(b(v_\omega))\,) = U_{v\omega}(\, g(v_\omega) \mid b_{v\omega}, s(b(v_\omega))\,). \tag{10}$$

By historical independence,

$$s(b(v_\omega)) = s(\, b(v_\omega) \mid g_{v\omega}). \tag{11}$$

Substitute the RHS of (11) into the RHS of (9) and then substitute the result into the LHS of (7). Substitute the LHS of (10) into the RHS of (7). Thus,

$$U_{v\omega}(\, b(v_\omega) \mid g_{v\omega}, s(\, b(v_\omega) \mid g_{v\omega})\,) > U_{v\omega}(\, b(v_\omega), s(b(v_\omega))\,). \tag{12}$$

For $b(v_\omega) \mid g_{v\omega}$, sequential rationality holds at ruler vertices succeeding $v_\omega$ since $v_\omega$ has no successors in Y, and historical independence assures that sequential rationality holds at $v_\omega$. So $b(v_\omega) \mid g_{v\omega}$ is sequentially rational. But since b is time consistent, it must be that $v_\omega$ is not in $R_1$ for otherwise b would be inconsistent at $v_\omega$. This argument holds for any last Y node, so no $v_\omega$ is in

$R_1$.

Claim(ii): No $v_\omega$ is in $R_2$. Again since $v_\omega$ is last, (12) holds and therefore b is sequentially irrational at $v_\omega$. Since we assume perfect policy formation and that b is nonidiosyncratic, $v_\omega$ must then be strategically relevant. No element of $R_2$ is strategically relevant, so no last Y node $v_\omega$ is in $R_2$.

Claim(i) and Claim(ii) establish that any last Y node $v_\omega$ must be in $R_3$. We now proceed by induction. Consider any $r \, \varepsilon \, R$ such that the only elements of Y to succeed it are in $R_3$.

Claim(iii): No r is in $R_1$. By supposition $b_r$ is a pure strategy. Let m be the citizen vertex $b_r$ leads to. Since r is in Y, $g_r$ must place weight on different (and better) moves at r. $g_r$ may be pure or it may distribute weight among equally good moves at r. At any rate, let c' be a generic element of the set of moves receiving positive weight from $g_r$, and let n be the citizen vertex c' leads to. For any vertex t, let NRS(t) be the set of *nearest ruler successor* vertices of t; excepting its endpoints, the curve connecting t and any element of NRS(t) contains no ruler vertices.[4]

We know a subgame originates at citizen node n since ruler actions are publicly observed. All the ruler vertices in G(n) are strategically irrelevant under b, so by nonidiosycrasy, b is sequentially rational in G(n). Let b" be any sequentially rational plan, let {b(NRS(n))} be the set of plans induced by b on the subgames originating with an element of NRS(n), and let { b"( NRS(r)\NRS(n)) } be the set of plans induced by b" on the subgames originating with an element of NRS(r) excepting the subgames originating with an element of NRS(n). Now we can construct a new plan:

$$h(r) := (\ g_r,\ \{\ b(NRS(n))\ \},\ \{\ b''(\ NRS(r) \backslash NRS(n)\ )\ \}\ )$$

Description of h(r): At r, h(r) specifies $g_r$; on subgames to which $g_r$ leads, h(r) specifies b; on subgames ensuing from r to which $g_r$ does not lead, h(r) specifies b''. Hence

$$U_r(\ h(r),\ s(b(r))\ ) = U_r(\ b(r) \mid g_r,\ s(b(r))\ ) \qquad (13)$$

By historical independence, $\{\ s(\ h(\ NRS(n)\ )\ )\ \} = \{\ s(\ b(\ NRS(n)\ )\ )\ \}$, which implies

$$U_r(\ h(r),\ s(b(r))\ ) = U_r(\ h(r),\ s(h(r))\ ) \qquad (14)$$

Since the only elements of Y to succeed r are in $R_3$,

$$U_r(\ b(r) \mid g_r,\ s(\ b(r)\ )\ ) = U_r(\ g(r),\ s(b(r))\ ) \qquad (15)$$

Substitute the RHS of (14) into the LHS of (13) and the RHS of (15) into the RHS of (13):

$$U_r(\ h(r),\ s(\ h(r)\ )\ ) = U_r(\ g(r),\ s(\ b(r)\ )\ ) \qquad (16)$$

Now, again since the only elements of Y to succeed r are in $R_3$,

$$U_r(\ g(r) \mid b_r,\ s(\ b(r)\ )\ ) = U_r(\ b(r),\ s(\ b(r)\ )\ ) \qquad (17)$$

Since r is an element of Y,

$$U_r(\ g(r),\ s(\ b(r)\ )\ ) > U_r(\ g(r) \mid b_r,\ s(\ b(r)\ )\ ) \qquad (18)$$

Substitute the LHS of (16) into the LHS of (18), and the RHS of (17) into the RHS of (18):

$$U_r(\ h(r),\ s(\ h(r)\ )\ )\ >\ U_r(\ b(r),\ s(\ b(r)\ )\ ) \qquad (19)$$

Now, $b(r)$ is sequentially rational at every ruler vertex in $G(r)$ except possibly $r$. Let

$$f^*_r\ :=\ \text{argmax over all } q_r \text{ of } U_r(\ h(r)\,|\,q_r,\ s(\ h(r)\ )\ ),$$

and let

$$f(r)\ :=\ h(r)\,|\,f^*_r.$$

$f(r)$ is sequentially rational.

Since $b_r$ is one option in the maximization problem defining $f^*_r$, (and since $s(f(r)) = s(h(r))$ by historical independence),

$$U_r(\ f(r),\ s(\ f(r)\ )\ )\ \geq\ U_r(\ h(r),\ s(\ h(r)\ )\ ) \qquad (20)$$

So by (19) and (20),

$$U_r(\ f(r),\ s(\ f(r)\ )\ )\ >\ U_r(\ b(r),\ s(b(r)\ )\ )$$

So if the players, following $(b, s(b))$, were to arrive at $r$, there the ruler would want to switch to the sequentially rational plan $f(r)$. Since $b$ is time consistent, there must be zero probability of $(b, s(b))$ reaching $r$, or $r$ is not in $R_1$.

<u>Claim(iv): r is not in $R_2$.</u> As in Claim(iii), since the only elements of Y to succeed $r$ are in $R_3$, equations (15) and (17) hold, and since $r$ is in Y inequality (18) holds. By historical independence $s(b(r)) = s(b(r)\,|\,g_r)$, so by substituting into the LHS of (15):

$$U_r(\ b(r)\,|\,g_r,\ s(\ b(r)\,|\,g_r)\ )\ =\ U_r(\ g(r),\ s(b(r)\ )\ ) \qquad (21)$$

Substitute the LHS of (21) into the LHS of (18), and substitute the RHS of (17) into the RHS of (18):

$$U_r(\ b(r)\,|\,g_r,\ s(\ b(r)\,|\,g_r)\ )\ >\ U_r(\ b(r),\ s(\ b(r)\ )\ ) \tag{22}$$

Thus r must not be in $R_2$, for otherwise r would be strategically irrelevant under b and (22) would contradict nonidiosyncrasy.

Review. Claim(i) showed that no last Y vertices are in $R_1$, and Claim(ii) showed that no last Y vertices are in $R_2$. Where r is a Y vertex such that the only Y vertices to succeed it are in $R_3$, Claim(iii) showed that r is not in $R_1$, and Claim(iv) showed that r is not in $R_2$. Thus by induction we conclude that all elements of Y are in $R_3$. Since no element of $R_3$ is in X, Y∩X is empty, and by Lemma 1 b is a best response to s(b). Hence (b, s(b)) is a Nash equilibrium.

◆◆◆

# NOTES

[1] Rasmusen, in his recent game theory text (1989, p82), notes that in an appropriately specified extensive game Stackelberg duopoly is a Nash equilibrium, but then suggests another equilibrium concept, much in line with this paper, to distinguish the spirit of Stackelberg duopoly from the spirit of Nash equilibrium: "An alternative definition is that a Stackelberg equilbrium is a strategy combination in which players select strategies in a given order [the ruler announces a plan!], and in which each player's strategy is a best response to the fixed strategies [re: plans] of the players preceding him and to the yet to be chosen strategies of players succeeding him, i.e., a situation in which players precommit to strategies in a given order. Such an equilibrium would not generally be either Nash or perfect."

[2] An note of interpretation: Suppose b is time inconsistent at v and $b_v$ is a mixed strategy. If the ruler reverts to an SR plan b' at v, and $b'_v$ specifies an action which received positive probability under $b_v$, how would the citizens know that the ruler changed her plan? We must assume that the citizens are fully informed of the new plan; they see the new roulette wheel that is spun at v, as well as the wheels to be spun at successor ruler nodes. The authors are grateful to Stergios Skaperdas for pointing this out.

[3] We are indebted to an anonymous referee for this example.

[4] Formally, $NRS(t) := \{v \in R: t < v$ and there does not exist any $v' \in R$ such that $t < v' < v\}$.

# REFERENCES

Blanchard, Olivier Jean and Stanley Fischer [1989]: *Lectures on Macroeconomics* (Cambridge, MA: MIT Press).

Fershtman, Chaim [1989]: "Fixed Rules and Decision Rules: Time Consistency and Subgame Perfection," *Economics Letters*, 30, 191-194.

Fischer, Stanley [1980]: "Dynamic Inconsistency, Cooperation and the Benevolent Dissembling Government," *Journal of Economic Dynamics and Control*, 2, 93-107.

Guiso, Luigi and Daniele Terlizzese [1990]: "Time Consistency and Subgame Perfection: The Difference between Promises and Threats," Banca d'Italia Discussion Paper, no. 138.

Hillier, Brian and James M. Malcomson [1984]: "Dynamic Inconsistency, Rational Expectations, and Optimal Government Policy," *Econometrica*, 52, 1437-1451.

Klein, Daniel and Brendan O'Flaherty [forthcoming]: "A Game-Theoretic Rendering of Promises and Threats," *Journal of Economic Behavior and Organization.*

_____ [1992]: "Imperfect Policy Formation and Time Inconsistency," Irvinve Economic Papers 90-91-03.

Kreps, David M. and Robert Wilson [1982]: "Sequential Equilibria," *Econometrica*, 50, 863-94.

Kydland, Finn and Edward Prescott [1977]: "Rules Rather Than Discretion: The Inconsistency of Optimal Plans," *Journal of Political Economy*, 85, 473-493.
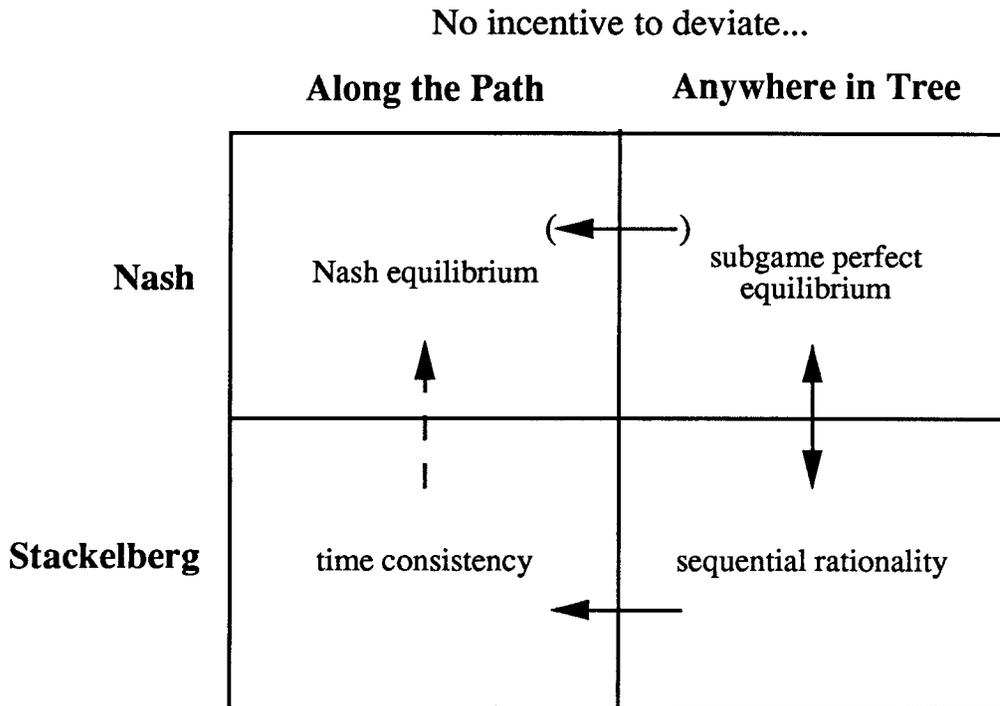
McTaggart, Douglas and David Salant [1988]; "Time Consistency and Subgame Perfect Equilibria in a Monetary Policy Game," ms.

Nalebuff, Barry and Martin Shubik [1988]: "Revenge and Rational Play," Princeton University discussion paper 138.

Rasmusen, Eric [1989]: *Games and Information* (Cambridge, MA: Basil Blackwell).

Schelling, Thomas C. [1960]: *The Strategy of Conflict* (Cambridge: Harvard University Press).
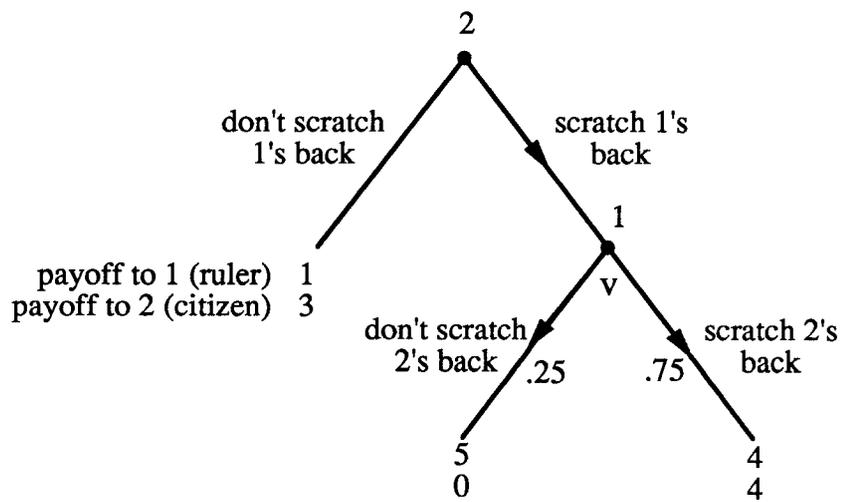
*Figure 1: Summary of Results*

No incentive to deviate...

|  | **Along the Path** | **Anywhere in Tree** |
|---|---|---|
| **Nash** | Nash equilibrium | subgame perfect equilibrium |
| **Stackelberg** | time consistency | sequential rationality |

The arrow in parantheses indicates a relationship that is well known.
This paper establishes the other four arrows. The dashed arrow signifies
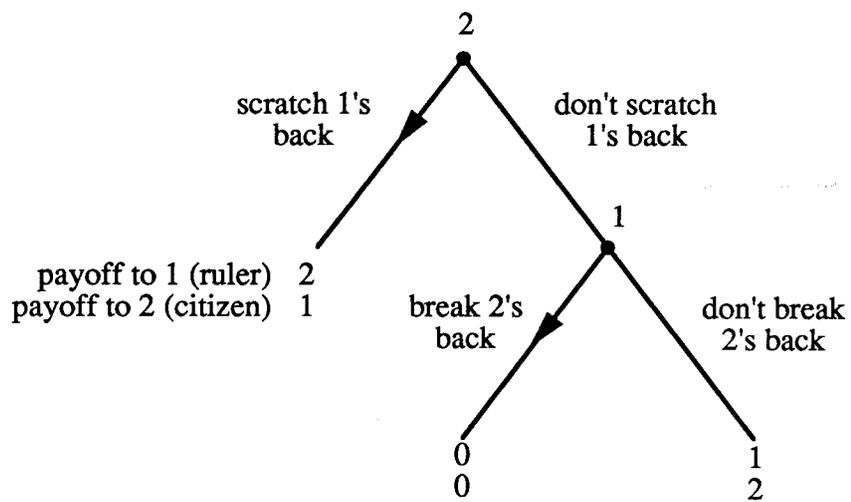that the relationship is true only with nontrivial restrictions.

## *Figure 2*

"If you scratch my back I'll scratch your back."

2

don't scratch
1's back

scratch 1's
back

1

payoff to 1 (ruler)   1
payoff to 2 (citizen)   3

v

don't scratch
2's back

.25

.75

scratch 2's
back

5
0

4
4

b and s(b) are shown by the arrows; b specifies (don't
scratch) with .25 probability and (scratch) with .75
probability. b is optimal and sequentially irrational, and it
is time inconsistent because at v 1 would like to switch to
sequentially rational play (don't scratch). (b, s(b)) is not
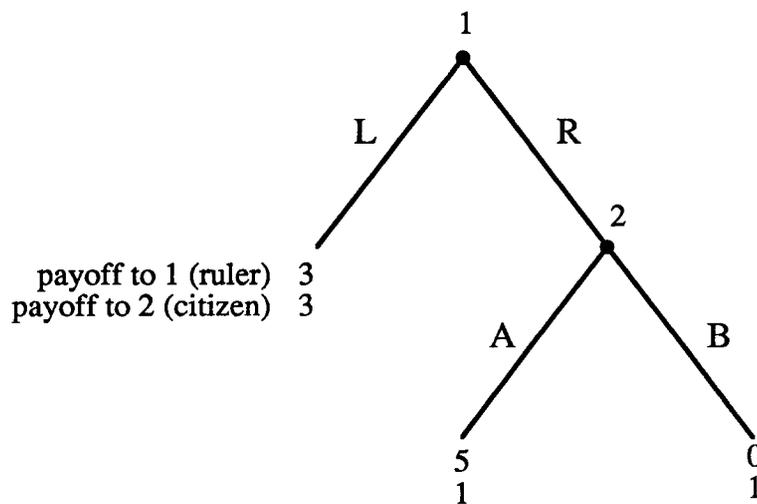a Nash equilibrium.

## *Figure 3*

"Scratch my back or else I'll break your back."



b and s(b) are shown by the arrows. b is optimal,
sequentially irrational, and time consistent. (b, s(b)) is a
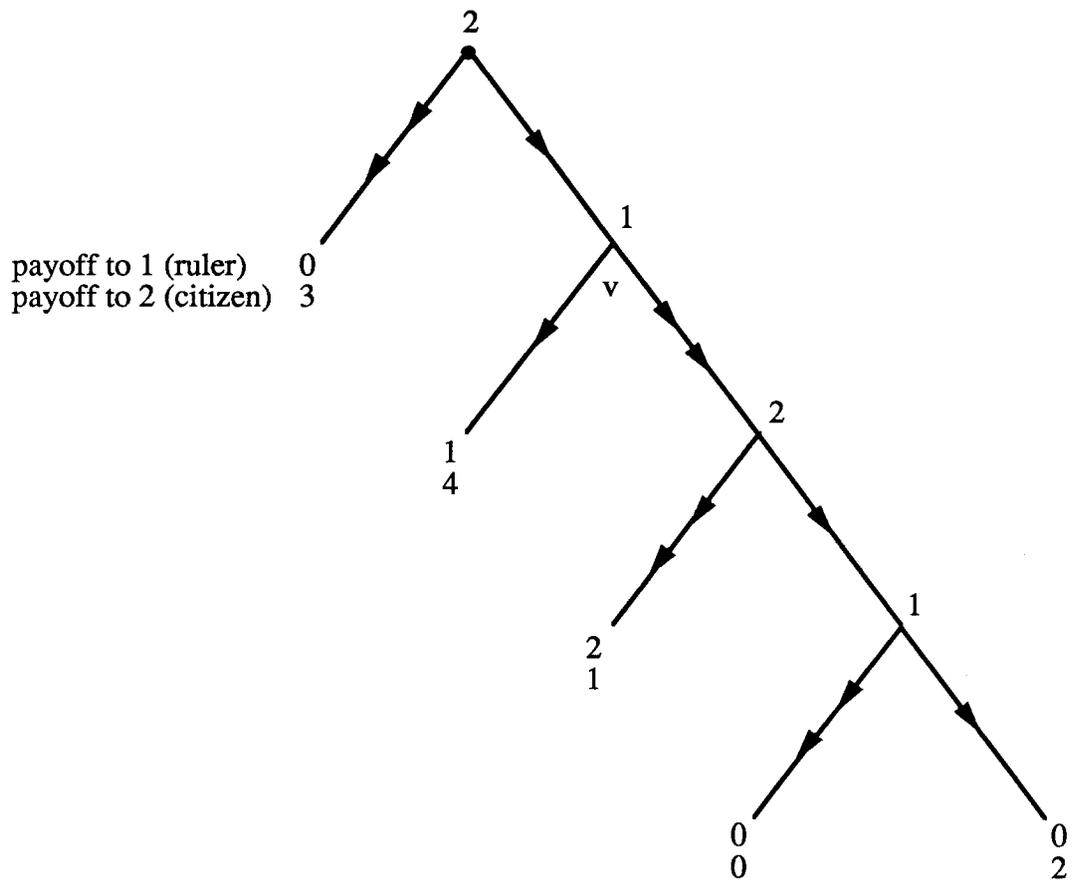Nash equilibrium.

## Figure 4

```
                              1
                              ●
                         L  /   \  R
                          /       \
                        /           2
payoff to 1 (ruler)  3 /          ●
payoff to 2 (citizen) 3      A  /   \  B
                            /       \
                          5           0
                          1           1
```

Violations of historical independence: Suppose the citizen's response is (B) whenever the ruler's plan places positive probability on (L), and (A) otherwise. Letting b = (L), (b, s(b)), that is, (L,B), is a subgame perfect equilibrium, but b is sequentially irrational.

Suppose the citizen's response is (B) whenever the ruler's plan places positive probability on (R), and (A) otherwise. Letting b = (L), b is sequentially rational, but (b, s(b)), that is, (L,A), is not a subgame perfect equilibrium.
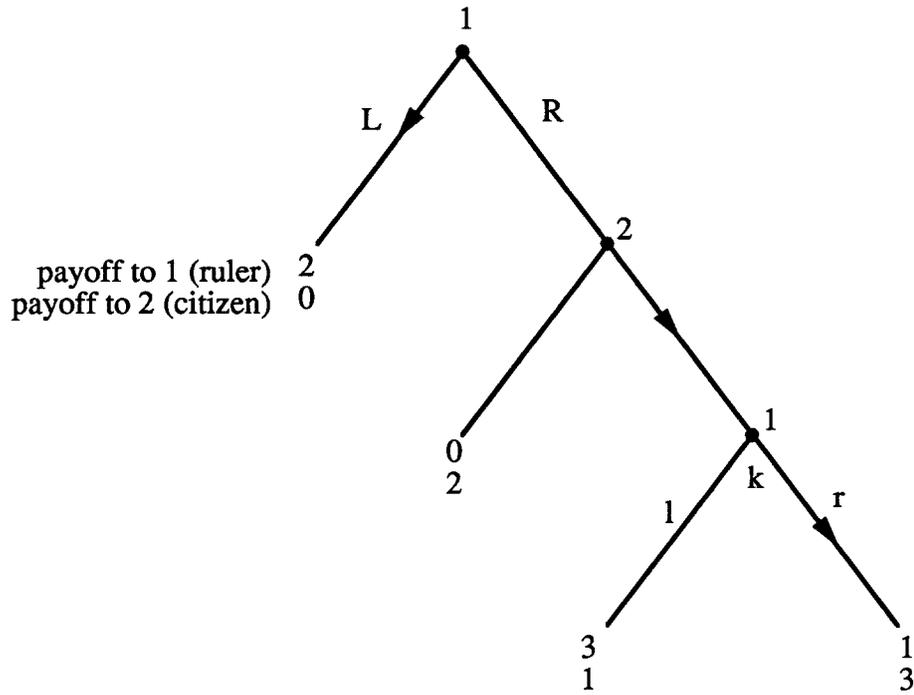
These public response rules violate historical independence and show that historical independence is necessary to Theorems 2 and 3.
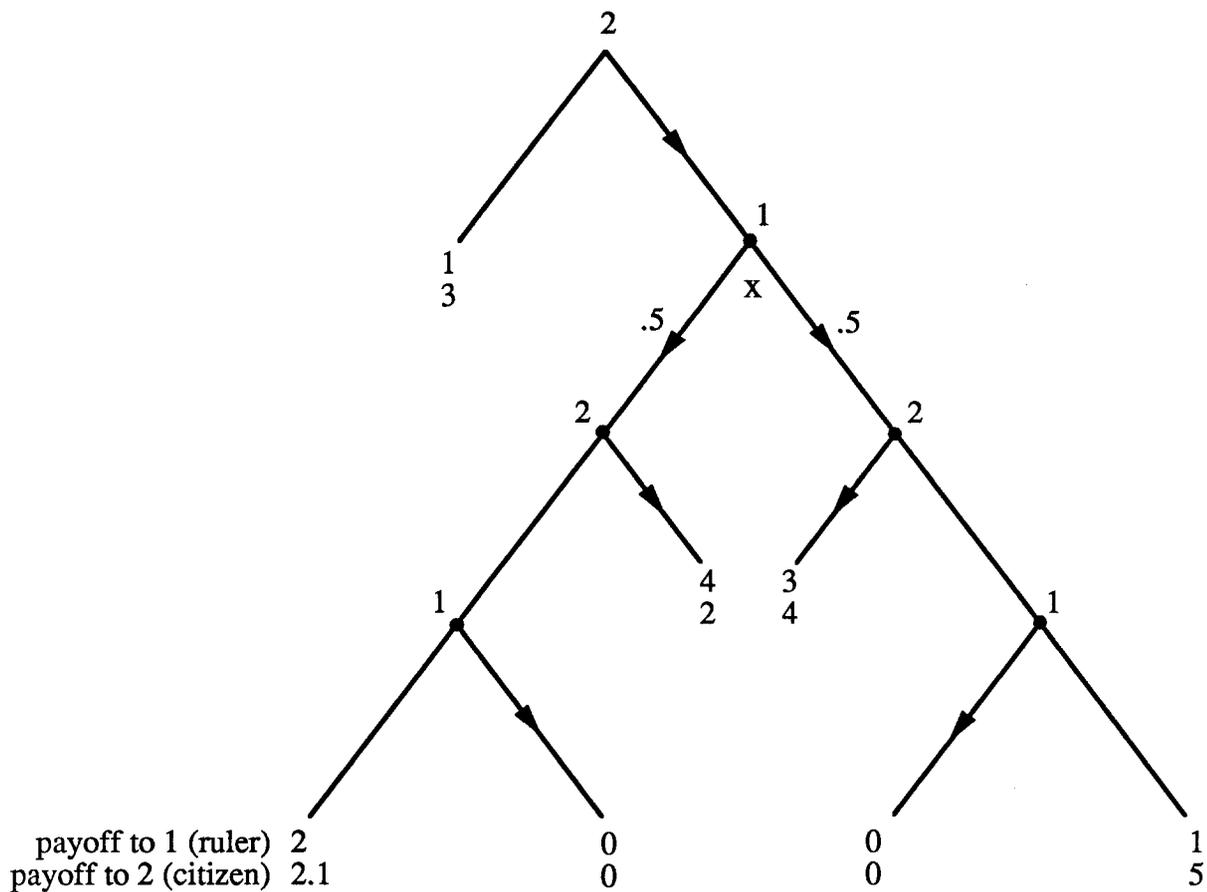
*Figure 5*

payoff to 1 (ruler)    0
payoff to 2 (citizen)   3

Violation of SR-equivalence:  b and s(b) are shown by the
single arrows.  b' and s(b') are shown by the double arrows.
Both b and b' are sequentially rational, yet $u_v(b(v)) = 1 < 2 = u_v(b'(v))$,
violating SR-equivalence.  b is time inconsistent since at v the ruler
would like to switch to the sequentially rational plan b'.  b is also
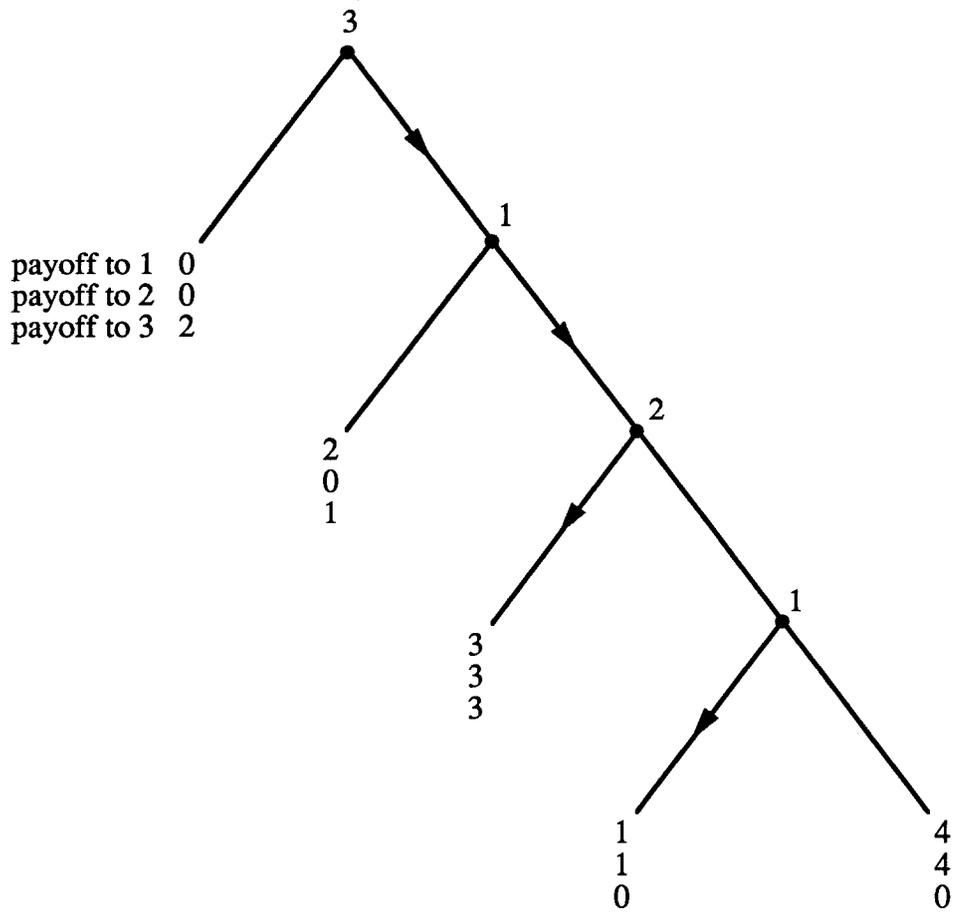optimal.

## Figure 6



b and s(b) are shown by the arrows.  b is optimal and time consistent.  (b, s(b)) is not a Nash equilibrium.  Node k is strategically irrelevant under b.  Since b is sequentially irrational at a strategically irrelevant node b is idiosyncratic.

## Figure 7



payoff to 1 (ruler)  2             0             0             1
payoff to 2 (citizen)  2.1          0            0             5

b and s(b) are shown by the arrows; at x b specifies each
action with .5 probability. b is optimal, time consistent and
nonidiosyncratic, yet (b, s(b)) is not a Nash equilibrium.
This example shows that the pure strategy condition in
Theorem 4 is necessary. (The point could have been
illustrated by the simpler game that results if you delete the
first move. The first move is included to show that
imposing that b be optimal, rather than that b be pure, will
not make Theorem 4 work.)

## Figure 8



payoff to 1   0
payoff to 2   0
payoff to 3   2

b and s(b) are shown by the arrows. b is optimal, nonidiosyncratic, and time inconsistent. (b, s(b)) is a Nash equilibrium.