

Sequential Equilibria

Author(s): David M. Kreps and Robert Wilson

Source: *Econometrica*, Vol. 50, No. 4 (Jul., 1982), pp. 863-894

Published by: [The Econometric Society](#)

Stable URL: <http://www.jstor.org/stable/1912767>

Accessed: 07/12/2010 13:12

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/action/showPublisher?publisherCode=econosoc>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



The Econometric Society is collaborating with JSTOR to digitize, preserve and extend access to *Econometrica*.

SEQUENTIAL EQUILIBRIA¹

BY DAVID M. KREPS AND ROBERT WILSON

We propose a new criterion for equilibria of extensive games, in the spirit of Selten's perfectness criteria. This criterion requires that players' strategies be *sequentially rational*: Every decision must be part of an optimal strategy for the remainder of the game. This entails specification of players' *beliefs* concerning how the game has evolved for each information set, including information sets off the equilibrium path. The properties of sequential equilibria are developed; in particular, we study the topological structure of the set of sequential equilibria. The connections with Selten's trembling-hand perfect equilibria are given.

1. INTRODUCTION

WE PROPOSE A NEW CRITERION for equilibrium in extensive games. The force of this criterion stems from the stringent requirement of *sequential rationality* imposed on the behavior of every player: Every decision must be part of an optimal strategy for the remainder of the game. In games with imperfect or incomplete information, this requirement entails conformity with Savage's [13] axioms of choice under uncertainty: At every juncture the player's subsequent strategy must be optimal with respect to some assessment of the probabilities of all uncertain events, including any preceding but unobserved choices made by other players. Mathematically, this is accomplished by broadening what is meant by an equilibrium. An equilibrium is not simply a strategy, but consists instead of two types of probability assessments by the players: the *beliefs* of a player concerning where in the game tree he is whenever it is his turn to choose an action, and his conjecture concerning what will happen in the future as given by the strategy. The novel aspect is the specification of beliefs on information sets that lie off the equilibrium path (that is, that have prior probability zero in the equilibrium). The specification of these beliefs allows us to verify that the player's own strategy is optimal starting from *every* point in the tree.

In this exposition we consider only games in which each player has perfect recall; cf. Kuhn [7]. For such games, sequential equilibria admit the following construction: In each player's personal decision tree induced by the game tree (cf. Wilson [16]) there is an appropriate assessment of the probabilities assigned

¹We have depended greatly and borrowed unrestrainedly from the insightful work of John Harsanyi and Reinhard Selten. We also gratefully acknowledge the influence of colleagues who have been exploring related ideas and who have discussed those ideas with us: Drew Fudenberg, Elon Kohlberg, Paul Milgrom, Roger Myerson, Roy Radner, John Roberts, Robert Rosenthal, Ariel Rubinstein, Sylvain Sorin, Jean Tirole. We are also grateful to two anonymous referees, who greatly aided the quality of exposition in the paper and pointed out an error in an earlier version of the paper. This research has been supported in part by National Science Foundation Grants SOC77-07741-A01 and SES80-06654 to the Institute for Mathematical Studies in the Social Sciences, Stanford University, by National Science Foundation Grant SES80-06407 to the Graduate School of Business, Stanford University, and by Office of Naval Research Grant No. 14-79-C-0685 to the IMSSS, Stanford University.

to *all* conditional uncertain events for which his strategy is among the optimal responses obtained by backwards recursion via dynamic programming. (These assessments are also required to be consistent among the players, in a fashion to be explained.) This encompasses the formulation of games with incomplete information due to Harsanyi [3]: A sequential equilibrium provides at each juncture an equilibrium in the subgame (of incomplete information) induced by restarting the game at that point.

Our definition of a sequential equilibrium recasts and slightly weakens Selten's [15] definition of a *perfect equilibrium*. Selten's definition accomplishes two things at once: It implicitly generates beliefs at information sets off the equilibrium path, and it requires that players' strategies be optimal with respect to those beliefs. In addition, it eliminates from consideration strategies that are otherwise weakly dominated. In a sequential equilibrium, the former is explicitly done, and the latter is dropped. Thus every perfect equilibrium is sequential, but not conversely. We prove, however, that if the former is done, then it is rarely the case that the latter is necessary; for "almost all" games the perfect and sequential equilibria "nearly" coincide. Thus, generically, the two concepts are identical mathematically.

We have two motives for proposing this alteration of Selten's definition. The first is pragmatic: In many examples of interest (e.g., in Kreps and Wilson [6], Milgrom and Roberts [8], and Rubinstein [12]), it is vastly easier to verify that a given equilibrium is sequential than that it is perfect. Second, making explicit the construction of beliefs off the equilibrium path enables discussion of which beliefs are "plausible" and which are not. Such discussion is difficult in the context of Selten's mechanical and indirect procedure for generating beliefs. And such comparisons can often help one to choose among sequential/perfect equilibria. (An example is given in Kreps and Wilson [6], where there is unique along-the-equilibrium-path behavior among all sequential equilibria whose beliefs meet an intuitively plausible monotonicity condition. Another example is given here in Section 8.) Indeed, we have found that by making the idea of beliefs explicit, the concept of a sequential equilibrium becomes consonant with the received tradition of single-person decision theory, and so it is easier to explain to nonspecialists.

The paper is organized as follows. In Section 2 a formulation of extensive games with perfect recall is given. The definitions of Nash equilibrium and Selten's concept of subgame-perfect equilibrium (Selten [14]) are reviewed in Section 3. In Section 4, we present examples that motivate restrictions more severe than subgame-perfection. Then we give the key definitions of *beliefs* and an *assessment*, and we formally define a *sequentially rational assessment*.

A sequential equilibrium is a sequentially rational assessment that meets certain further consistency criteria. For example, beliefs along the equilibrium path should be computed from the strategy via Bayes' rule. In Section 5 we present the consistency criterion that we subsequently use, and we motivate this criterion by a series of examples.

Properties of sequential equilibria are presented in Section 6. Basic properties are established: They are Nash, exist for all games, and have an upper hemi-continuous correspondence over the space of payoffs. Then we consider their relation to Selten's subgame-perfection concept and an extension of this concept. Lastly, we study the structure of the set of sequential equilibria for "generic" games.

In Section 7 we compare Selten's criterion for (trembling-hand) perfect equilibrium with our criterion for sequential equilibrium. The basic results are as stated above: Every perfect equilibrium is sequential. Conversely, for generic payoffs the sets of perfect and sequential equilibria "nearly" coincide. (The possible exceptions are limited to weak equilibria; and there is complete coincidence of the sets of "equilibrium paths.") This generic equality of sequential and perfect equilibria becomes an exact identity *if* one weakens the defining apparatus of perfect equilibria to allow perturbations of the payoffs.

In Section 8 we develop the idea that explicit consideration of the beliefs off the equilibrium path can help one to choose among sequential equilibria. Concluding remarks are made in Section 9.

The technical results of Sections 6 and 7 are proved in an Appendix. The Appendix also gives several additional characterizations that are not discussed in the text.

2. EXTENSIVE GAMES

We will use a formulation of an extensive game that is equivalent to Kuhn's [7]. In this formulation, the following are specified: (1) the physical order of play; (2) the choices available to a player whenever it is his turn to move; (3) rules for determining whose move it is at any point; (4) the information a player has whenever it is his turn to move; (5) the payoffs to the players as functions of the moves they select; (6) the initial conditions that begin the game (that is, the actions of nature).

We illustrate our formulation with the following example of a game with two players. Player 1 moves first and chooses between three actions: L, R, A . If player 1 chooses A , then the game ends, with payoffs zero to each player. If 1 chooses L or R , then it is 2's turn to choose between actions l and r , after which the game ends, and payoffs are made. When and if 2 does get to choose between l and r , 2 does not know which of L and R 1 chose—2 knows only that 1 did not choose A . A diagrammatic representation of this is given in Figure 1.

Mathematically, the formulation is constructed from the following objects:

(1) The physical order of play is given by a finite set T of *nodes* together with a binary relation $<$ on T that represents *precedence*. In the example, the set T consists of eight points: the open circle, the two closed circles, and the four column vectors. Precedence is indicated by arrows—one node precedes another if there is a sequence of arrows pointing from the first to the second. The binary relation $<$ must be a partial order, and $(T, <)$ must form an arborescence: The

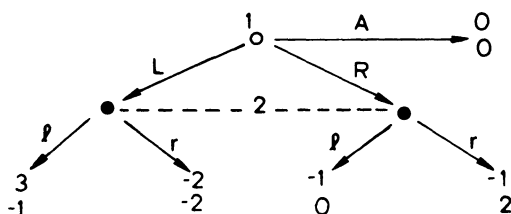


FIGURE 1.

relation $<$ totally orders the predecessors of each member of T . (This prevents cycles from appearing in the order of play, and it means that each node in the tree can be reached by one and only one path from an initial node through the tree.) Useful auxiliary notation and definitions are compiled in Table I.

TABLE I
NOTATION AND DEFINITIONS FOR EXTENSIVE GAMES

Name	Notation	Definition
terminal nodes or outcomes	Z	$\{t \in T : t \text{ has no successors}\}$
decision nodes	X	$T \setminus Z$
initial nodes or states	W	$\{t \in T : t \text{ has no predecessors}\}$
predecessors of t	$P(t)$	$\{x \in X : x < t\}$
immediate predecessor of t	$p_1(t)$	$\max\{x : x < t\}$ for $t \notin W$
n th predecessor of t	$p_n(t)$	$p_1(p_{n-1}(t))$ for t such that $p_{n-1}(t) \notin W$; $p_0(t) = t$ for all t
number of predecessors of t	$l(t)$	$l(t)$ is such that $p_{l(t)}(t) \in W$
immediate successors of x	$S(x)$	$p_1^{-1}(x)$ for $x \in X$
terminal successors of x	$Z(x)$	$\{z \in Z : x < z\}$ for $x \in X$

Note that we depict the terminal nodes (nodes in Z) by column vectors (the meaning of which will be discussed below), the decision nodes (nodes in X) by circles, and the initial nodes (nodes in W) by open circles. In the example there is a single initial node, but in general there may be more than one element of W —see the game depicted in Figure 8 for an example.

The following interpretations are made: The game begins at one of the initial nodes (determined by nature—see below) and then proceeds along some path from node to immediate successor, terminating when a terminal node is reached. The various paths give the various possible orders of play. In many games of interest, especially games with simultaneous moves, more than one tree can be used to represent the game.

(2) To represent the choices available to players at decision nodes, we have a finite set A of actions and a function $\alpha : T \setminus W \rightarrow A$ that labels each non-initial node with the last action taken to reach it. In the figures, the actions are labelled along the branches, so that the L on the uppermost left-hand branch is read as saying that action L leads from the initial node to the node below and to the left of it. In terms of the function α , L is the value of α at the node below and to the left of the initial node. Note that $\alpha(S(x))$ is thus the set of feasible actions at the

decision node x . For example, the set of actions feasible at the initial node is $\{L, R, A\}$. We require that α be one-to-one on the set $S(x)$ of immediate successors of x .

(3) To represent the rules for determining whose move it is at a decision node, we have a finite set I of *players* and a function $\iota: X \rightarrow I$ that assigns to each decision node the player whose turn it is. In Figure 1, the 1 above the initial node records that it is player 1's move at that point, and the 2 connected to the two other decision nodes by the dashed line records that it is 2's move at each of those two points.

(4) Information possessed by players is represented by a partition H of X that divides the decision nodes into *information sets*. The cell $H(x)$ of H that contains x identifies the decision nodes that player $\iota(x)$ cannot distinguish from x based on the information he has available when it is his turn to choose an action at x . We depict information sets by connecting nodes in a single information set with a dashed line. Thus in Figure 1, the dashed line connecting the two non-initial decision nodes denotes that these two nodes lie in the same information set—player 2, when it is his turn to choose an action, doesn't know whether 1 chose L or R . We require that a player knows when it is his turn to choose and which actions are feasible:

$$(2.1) \quad \text{If } x \in H(x'), \text{ then } \iota(x) = \iota(x') \text{ and } \alpha(S(x)) = \alpha(S(x')).$$

Thus it makes sense to write $\iota(h)$ and to partition H into sets $H^i = \iota^{-1}(i)$. (In our figures, we label information sets instead of nodes with players.) It also makes sense to write $A(h)$ for $\alpha(S(h))$, the set of actions feasible at information set h . (Note that 2's feasible actions at the two non-initial nodes are identical. What we really have is that 2 can choose between l and r at his *information set*.) For notational convenience, we assume that α is onto and that for each $a \in A$, $A^{-1}(a)$ is a singleton set. That is, each action can be taken only in a single information set. Then we can partition A into sets $A^i = \{a: A^{-1}(a) \subseteq H^i\}$ for $i \in I$.

We also assume that each player has perfect recall. Each player knows whether he chose previously:

$$(2.2) \quad \text{If } x \in H(x'), \text{ then } x \not\prec x'.$$

And he knows whatever he knew previously, including his previous actions:

$$(2.3) \quad \text{If } x, x', x'' \in \iota^{-1}(i), x \prec x', \text{ and } H(x') = H(x''), \text{ then } H(x)$$

includes some predecessor of x'' at which the same action was chosen

as was chosen at x ; that is $P(x'') \cap H(x) = \{x^0\}$, and if $x = p_n(x')$

and $x^0 = p_m(x'')$, then $\alpha(p_{n-1}(x')) = \alpha(p_{m-1}(x''))$.

The collection $\{T, <; A, \alpha; I, \iota; H\}$ defines an *extensive form*. To obtain an *extensive game* we add a specification of the players' utilities assigned to the terminal nodes and the probabilities assigned to the initial nodes.

(5) For each player i , the *payoff function* $u^i : Z \rightarrow R$ assigns a real-valued von Neumann-Morgenstern utility to each outcome. We denote a specification of the payoffs by $u = (u^i(z)) \in R^{I \times Z}$. In our pictorial representations, we simply write the vector of payoffs to the players (player 1 first, player 2 second, etc.) for the terminal node. For example, in Figure 1, if player 1 plays L and 2 plays l , then the payoffs are 3 to player 1 and -1 to player 2.

(6) Player i 's *initial assessment* ρ^i is a probability measure on the set W of states or initial nodes. (For notational convenience, we have put all actions by nature at the "start" of the game.) To keep matters simple, we henceforth assume that the players' initial assessments are strictly positive and are all the same: $\rho^i \equiv \rho \gg 0$. When necessary, we depict initial assessments by recording the probability $\rho(w)$ in braces next to the node w .

A *pure strategy* for player i is an assignment $\sigma^i : H^i \rightarrow A$ such that $\sigma^i(h) \in A(h)$. This specifies what action player i will take each time it is his turn to choose, based on the information that he possesses. One defines a *mixed strategy* for player i as a probability distribution over the set of his pure strategies. However Kuhn [7] shows that for games with perfect recall it is sufficient to restrict attention to *behavior strategies*, hereafter simply called strategies.

(7) A *strategy* $\pi^i : A^i \rightarrow [0, 1]$ for player i assigns to each information set $h \in H^i$ a probability measure on the set $A(h)$. That is, $\sum_{a \in A(h)} \pi^i(a) = 1$, for each $h \in H^i$.

Let Π^i denote the set of strategies for player i , and let $\Pi = \times_{i \in I} \Pi^i$ be the set of *strategies* for the game. Each strategy $\pi \in \Pi$ induces a probability measure P^π on the set Z of outcomes according to the formula:

$$P^\pi(z) = \rho(p_{l(z)}(z)) \prod_{l=1}^{l(z)} \pi^{i(p_l(z))}(\alpha(p_{l-1}(z))).$$

The expectation operator using P^π is denoted $E^\pi[\cdot]$; in particular, $E^\pi[u^i(z)]$ is player i 's expected utility from the strategy π . We shall frequently use the notation $P^\pi(x)$ and $P^\pi(h)$ for $P^\pi(Z(x))$ and $P^\pi(Z(h))$, respectively.

A *subform* of an extensive form is a collection of nodes $\hat{T} \subseteq T$, together with $<, \iota, A, \alpha$ and H all defined on \hat{T} by restriction, satisfying closure under succession and preservation of information sets: $S(x) \subseteq \hat{T}$ and $H(x) \subseteq \hat{T}$ if $x \in \hat{T}$. For every node $x \in X$ there is a minimal subform $\hat{T}(x)$ containing x . Note that x need not be an initial node in $\hat{T}(x)$ (cf. Figure 8). A *proper subform* (following Selten [14]) is a subform \hat{T} consisting solely of some node x and its successors. In this case we call x the *root* of \hat{T} . Given a proper subform \hat{T} with root x , there is a well-defined proper subgame starting with x as the unique initial node. That is, the game is formed by \hat{T} and all the structure that \hat{T} inherits from the original form, the payoffs u restricted to $\hat{T} \cap Z$, and the initial assessment $\hat{\rho}(x) = 1$. For nonproper subforms a subgame is not always well-

defined, if the initial assessment $\hat{\rho}$ is lacking. (Example: in Figure 1, the two nodes that form player 2's information set together with the four terminal nodes that follow constitute a nonproper subform. Note that if we tried to define a subgame for this subform, we would be lacking an initial assessment on the two (now initial) nodes.)

3. NASH EQUILIBRIA AND SUBGAME PERFECTION

The weakest criterion for equilibrium that we shall discuss is the familiar one proposed originally by Nash [11]. A strategy is a *Nash equilibrium* if each player's strategy is an optimal response to the other players' strategies. That is, $\pi \in \Pi$ is a Nash equilibrium if, for each player $i \in I$,

$$E^\pi[u^i(z)] \geq E^{\bar{\pi}}[u^i(z)] \quad \text{for every strategy } \bar{\pi} \in \Pi$$

such that $\bar{\pi}^j = \pi^j$ for $j \neq i$.

This definition has been motivated in many ways, and we shall not attempt to repeat those motivations here. But a thread common to all of them is that if players are to arrive at some "agreed-upon" mode of behavior, then it is *necessary* that this behavior constitutes a Nash equilibrium. Otherwise, some player would find it advantageous to defect from the agreement. (The different interpretations vary in their explanations of how it might be that such an agreement would arise, whether such an agreement must be explicit, etc.)

Consideration of games in extensive form lead to other, more stringent necessary conditions for "agreed-upon" behavior. One such condition is Selten's [14] criterion of *subgame-perfection*. Consider the game depicted in Figure 2. One Nash equilibrium for this game has player 1 choosing *L* and player 2 choosing *l*. Note well that 1 chooses *L* because he anticipates the choice of *l* by 2, while 2 is content to choose *l* only because his choice is irrelevant so long as 1 chooses *L*. But if 1 were to choose *R*, then it seems reasonable to suppose that 2, facing this *fait accompli*, would choose *r*. And 1, realizing this, "should" choose *R*. Selten [14] has formalized this intuition in the following criterion.

DEFINITION: Strategy π is *subgame perfect* if for every proper subgame the strategy π restricted to the subgame constitutes a Nash equilibrium for the subgame.

This definition makes sense because in any proper subgame it makes sense to speak of each player's expected utility in that subgame, and thus the Nash

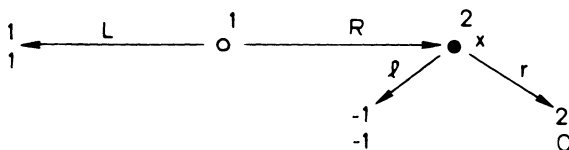


FIGURE 2.

criterion can be applied (in Figure 2, to the subgame with root x). (In Selten [14], nature's moves are not all put at the beginning of the tree T . This gives him "more" proper subgames and, correspondingly, more applications for the criterion above. The difference is insignificant. We have put nature's moves at the beginning of the tree for convenience only—nothing in the analysis changes if this is relaxed. We could also remedy matters by calling a subform \hat{T} proper if there is a unique probability distribution $\hat{\rho}$ on \hat{W} such that $\hat{\rho}(w)P^\pi[\hat{W}] = P^\pi[\hat{w}]$ for all $\hat{w} \in \hat{W}$ —then this $\hat{\rho}$ is the natural candidate for the initial assessment.) And it is a natural restriction for any "agreed-upon" behavior—otherwise the agreement would not hold up if the subgame were reached. Accordingly some player might defect from the agreement and cause the subgame to be reached, anticipating a breakdown of the agreement favorable to himself.

4. BELIEFS AND SEQUENTIAL RATIONALITY

Selten has gone on to observe that the intuitive motivation for the subgame perfection criterion can be applied to games that lack proper subgames. This is illustrated by the game depicted in Figure 3, taken from [15, Section 6]. One Nash equilibrium for this game has player 1 choosing D , player 2 choosing a , and player 3 choosing l . This equilibrium is subgame perfect, as the only proper subgame here is the game itself. But, as Selten argues, this equilibrium is not sensible. The behavior of player 2 is hard to justify, if it is supposed that 3 will choose l . Note that player 2's information set is a singleton, and therefore there is no difficulty in taking the other players' strategies as given and asking: If this node is reached, then what action is optimal for 2? That is, the conditional expected payoff to 2 on reaching x is calculable from the strategies of the other players. Given the supposed behavior of 3, 2 prefers to choose d . (The reader can verify that if 1 realizes this, then 1 would optimally choose A instead of D , thereby upsetting the equilibrium. The only "sensible" equilibrium in this game has 1 choosing A , 2 choosing a , and 3 choosing r with probability at least $3/4$.)

The subgame perfection criterion, as formally defined, fails in this example, because there is not a proper subgame starting from the node x —player 3 is unable to compute his expected utility in this subgame. But because this information set for 2 is a singleton, we can compute expected utility for 2

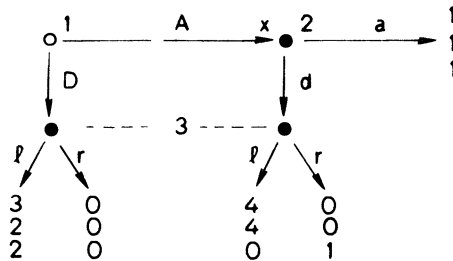


FIGURE 3.

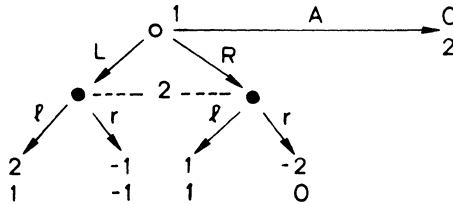


FIGURE 4.

conditional on hitting this information set, and this is enough to reject the supposed equilibrium. A corresponding general criterion can be formulated as follows: A strategy π should be such that for any information set h that is a singleton, player $u(h)$ should not be able to change his strategy unilaterally and thereby improve this expected utility starting from h .

The restriction of this criterion to singleton information sets h is necessary mathematically, so that player $u(h)$'s expected utility starting from h can be calculated. But it does limit the applicability of the criterion. Consider the game depicted in Figure 4. A Nash equilibrium for this game has 1 choosing A , and 2 choosing r . This strategy is subgame perfect, and it satisfies the further criterion given above. But if 1 gives the move to 2, then regardless of what 2 thinks the chances are that he is at one node or the other in his information set, 2 will do better by choosing l . And if 1 realizes this and concludes that 2 will choose l , then 1 will optimally pick L .

We are unable to apply the subgame perfection criterion or the other criterion above for technical reasons: The strategy π does not provide sufficient information to compute player 2's expected payoff conditional on reaching his information set. But if 2 is rational in the sense of Savage [13]—when faced with a choice 2 makes some assessment about what 1 did that is consistent with what 2 knows, and then optimizes accordingly—then 2 will choose l . This is the substance of sequential rationality: The strategy of each player starting from each information set must be optimal starting from there according to some assessment over the nodes in the information set and the strategies of everyone else.

To formalize this, as part of the description of an equilibrium we specify for each information set h the assessment made by player $u(h)$ over the nodes in h if h is reached. A *system of beliefs* is defined as a function $\mu: X \rightarrow [0, 1]$ such that $\sum_{x \in h} \mu(x) = 1$ for each $h \in H$. Interpret $\mu(x)$ as the probability assigned by $u(h)$ to $x \in h$ if h is reached. An *assessment* is a pair (μ, π) consisting of a system of beliefs μ and a strategy π . Given an assessment (μ, π) , for each $h \in H$ we can define "conditional" probability $P^{\mu, \pi}(\cdot | h)$ over Z in the obvious fashion:

$$\text{If } z \notin Z(h), \text{ then } P^{\mu, \pi}(z | h) = 0.$$

$$\text{If } z \in Z(h), \text{ say } p_n(z) \in h,$$

$$\text{then } P^{\mu, \pi}(z | h) = \mu(p_n(z)) \cdot \prod_{m=1}^n \pi(\alpha(p_{m-1}(z))).$$

(We shall use notation such as $P^{\mu,\pi}(h' | h)$ as shorthand for $P^{\mu,\pi}(Z(h') | h)$ when convenient.) Denoting conditional expectations by $E^{\mu,\pi}[\cdot | h]$, we say that the assessment is *sequentially rational* if, for all $h \in H$,

$$E^{\mu,\pi}[u^{i(h)}(z) | h] \geq E^{\mu,\bar{\pi}}[u^{i(h)}(z) | h]$$

for all $\bar{\pi}$ such that $\bar{\pi}^j = \pi^j$ for $j \neq i(h)$.

In words, taking the beliefs as fixed, no player prefers at any point to change his part of the strategy π .

A sequential equilibrium, roughly speaking, is a sequentially rational assessment (μ, π) . This is only a rough definition because we first want to impose consistency conditions, such as that assessments obey Bayes' rule when it applies. We develop these conditions and give the corresponding definition of a sequential equilibrium in the next section. The point to be stressed here is that an *assessment* (and not simply a strategy) will or will not be a sequential equilibrium. We require that an equilibrium specify beliefs as well as strategies.

Selten [15] suggests a somewhat different "cure" for the problem posed by examples such as the game depicted in Figure 4. He proposes a criterion called (among game theorists) "trembling-hand" perfection. We describe this criterion and its relation to ours in Section 7.

5. CONSISTENT ASSESSMENTS AND SEQUENTIAL EQUILIBRIA

We begin by giving the formal definition of a consistent assessment (μ, π) and the corresponding definition of a sequential equilibrium. Let Π^0 be the set of all strictly positive strategies. That is, $\pi \in \Pi^0$ if $\pi(a) > 0$ for all $a \in A$. If $\pi \in \Pi^0$, then $P^\pi(x) > 0$ for all x , and the only reasonable way to define beliefs μ associated with π is via Bayes' rule:

$$\mu(x) = P^\pi(x) / P^\pi(H(x)).$$

Let Ψ^0 denote that subset of the set of assessments (μ, π) where $\pi \in \Pi^0$ and μ is defined from ρ and π by Bayes' rule.

DEFINITION: An assessment (μ, π) is *consistent* if $(\mu, \pi) = \lim_{n \rightarrow \infty} (\mu_n, \pi_n)$ for some sequence $\{(\mu_n, \pi_n)\} \subseteq \Psi^0$. The set of consistent assessments is denoted by Ψ . (That is, Ψ is the closure of Ψ^0 .)

A *sequential equilibrium* is an assessment (μ, π) that is both consistent and sequentially rational.

This definition of consistency is not completely intuitive on its own. We propose it because it neatly embodies a number of distinct intuitive notions of consistency. To provide motivation, we now survey those more intuitive notions.

An obvious consistency criterion for any assessment (μ, π) is that μ must be

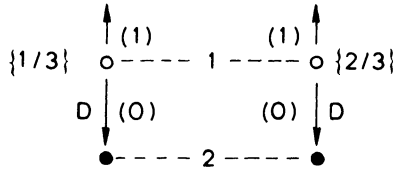


FIGURE 5.

defined from π by Bayes' rule whenever possible:

$$(5.1) \quad \mu(x)P^\pi(H(x)) = P^\pi(x).$$

Note that, since ρ is used to define P^π , consistency with ρ is embodied in this definition. This very basic criterion is clearly implied by $(\mu, \pi) \in \Psi$.

This uniquely defines $\mu(x)$ for any x such that $P^\pi(H(x)) > 0$. What happens when a player reaches an information set h with $P^\pi(h) = 0$? It is plausible to suppose that the player will construct some hypothesis as to how the game has been played, in the form of a strategy π' that satisfies $P^{\pi'}(h) > 0$ and then use π' and Bayes' rule to compute $\mu(x)$ for $x \in h$. This procedure limits the possible beliefs of a player. For example, in the part of a game depicted in Figure 5, player 2's beliefs in his information set must attach probability 1/3 to the left-hand node. This is because player 1 cannot distinguish between the two nodes in his information set, so any strategy he could hypothesize that gives 2's information set positive probability must (by Bayes' rule) preserve the initial probability assessment. That is, simply assuming that the players' beliefs always respect the informational structure of the game constrains players' beliefs. Formally:

$$(5.2) \quad (\mu, \pi) \text{ is structurally consistent if for each } h \in H \text{ there exists some strategy } \pi' \in \Pi \text{ such that } P^{\pi'}(h) > 0 \text{ and } \mu(x) = P^{\pi'}(x)/P^{\pi'}(h) \text{ for all } x \in h.$$

If $(\mu, \pi) \in \Pi$, then (μ, π) is structurally consistent. (A direct proof is easy.)

One can carry this "alternative hypothesis" story a step further. Fix a player i . His "primary hypothesis" as to how the game will be played is π , and if his beliefs obey (5.1), then he applies π to compute μ whenever possible. We might assume that when π does not apply—when he comes to an information set h with $P^\pi(h) = 0$ —then he has a "second most likely hypothesis" $\pi(2)$ that he attempts to apply. If that fails, he tries his "third most likely hypothesis" $\pi(3)$, and so on. Formally, we suppose that each player i has a finite sequence of hypotheses $\pi(1) = \pi, \pi(2), \pi(3), \dots, \pi(K)$, where for each $h \in H^i, P^{\pi(k)}(h) > 0$ for some $k \leq K$, and that $\mu(x)$ for $x \in h$ is computed using Bayes' rule applied to that $\pi(k)$ of lowest index k that satisfies $P^{\pi(k)}(h) > 0$. The force in this is that the sequence of "alternative hypotheses" is independent of h — $\pi(2)$ is the player's second most likely hypothesis for all $h \in H^i$. A further strengthening of this

requires that all players use the same finite sequence $\pi(1) = \pi, \pi(2), \dots, \pi(K)$. This requirement is in the spirit of the “common knowledge” hypothesis of Nash equilibrium—if there are rational secondary hypotheses, they should be unanimously held, just as is the primary hypothesis π . We call this strengthened consistency criterion *lexicographic consistency*.

A comparison of lexicographic consistency and our original definition of consistency is made easy by the following result. Let Δ be the set of all probability measures on Z of the form P^π for $\pi \in \Pi$.

LEMMA 1: *A sufficient condition for (μ, π) to satisfy lexicographic consistency is that there exists a sequence of probability measures $\{P_n\} \subseteq \Delta$ such that $\lim_n P_n = P^\pi$ and, for each x , $\mu(x) = \lim_n P_n(x) / P_n(H(x))$.*

The proof is left to the reader. (The methods used in the Appendix to prove Lemma A2 are easily adapted to this case.) The criterion embodied in this lemma is a bit stronger than lexicographic consistency, and the analogous criterion that is equivalent to lexicographic consistency is a bit cumbersome. (Essentially, one must allow P_n that are in the convex hull of Δ , and that are asymptotically “close” to Δ .) But since this criterion is clearly implied by our original definition of consistency, we see that lexicographic consistency is subsumed by that definition.

Consideration of some examples motivates further restrictions. Consider, for example, the part of an extensive game (with strategies and beliefs) that is depicted in Figure 6. (Beliefs are depicted in square brackets, and strategies in parentheses.) In particular, compare player 2’s beliefs in his two information sets. We claim that these beliefs are inconsistent with each other and with player 3’s strategy. For *if* player 2 reaches his first information set and adopts the beliefs shown, *then* he expects (given 3’s strategy) to reach his second information set. And if he uses Bayes’ rule starting from his first information set together with 3’s strategy in order to obtain his beliefs in the second, he would not come up with the beliefs shown. This does not violate Bayesian or lexicographic consistency; for the latter, the “secondary hypothesis” $\pi(2)$ simply has *both* players 1 and 3 changing their strategy. In this instance, what is wanted is an extension of Bayes’ rule. For a single player, this might be formulated as:

$$(5.3) \quad P^{\mu, \pi}(x' | h) = \mu(x') P^{\mu, \pi}(h' | h)$$

for all h, h' , and x' such that $u(h) = u(h')$, $h < h'$, and $x' \in h'$.

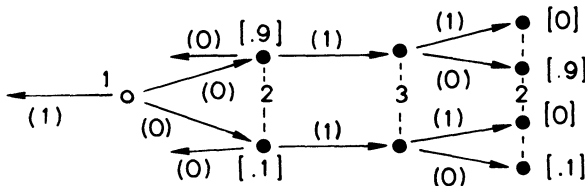


FIGURE 6.

In words, starting from *any* information set, a player uses his beliefs at that point together with π and Bayes' rule to compute subsequent beliefs when possible. The key here is the *continued* use of π after an initial defection. The philosophy behind this is that the strategy π at an information set h should encode the players' conjectures concerning what will happen *if* h is reached. *If* h is reached only by a defection, then π at h should encode what will subsequently happen, conditional on that initial defection. A first defection does not make a second more likely; correlation in defections are (partially) ruled out.

A referee has suggested an extension to (5.3), following the principle that rational beliefs should be common knowledge and thus commonly shared. If one accepts this, then in the spirit of (5.3) we could remove the restriction that $u(h) = u(h')$ as follows:

$$(5.4) \quad P^{\mu, \pi}(x' | h) \mu(x) = P^{\mu, \pi}(x | h) \mu(x')$$

for all h, h', x , and x' such that $x, x' \in h'$ and $x, x' \in S(h)$.

(The motivation for this will become clear if the reader draws the picture entailed and re-expresses the equation in ratio form.) One special case implied by (5.4) deserves mention. This is where h in (5.4) is a singleton set. For this case (5.4) can be paraphrased: Players use Bayes' rule applied to π in any proper subgames that arise. Of course, (5.4) implies (5.3) (sum over $x \in h'$), and both are implied by our general consistency condition.

Consider lastly the piece of the extensive game given in Figure 7. Player 3's beliefs are explicable as follows. Upon unexpectedly reaching his information set, he reconstructs the play of the game as follows: 1 changed his strategy to .9 for the upper branch, and 2 changed his to give positive probability to a move to the right. With this as $\pi(2)$, these beliefs are lexicographically consistent. And it is easy to verify that (5.4) is satisfied. But are 3's beliefs reasonable? If one grants the principle that defections from the equilibrium strategy ought to be uncorrelated—that given that 2 has defected (which he surely must have, given that 3's information set has been reached), the most likely (in a lexicographic sense) hypothesis is that 1 continues to play according to π —then the answer is no. Player 3 should give more credence to the hypothesis that only player 2 defected, and he should therefore have beliefs .1 at the top-most node. (This example becomes even more stark if we suppose that players 2 and 3 are the same. Then the defection by 2 is *known* to 3—should this player revise his assessments as to what 1 did because he himself defected?)

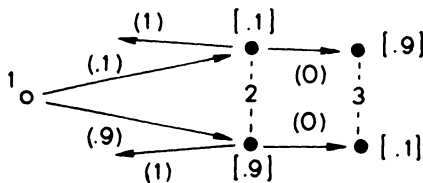


FIGURE 7.

Our initial consistency criterion is advanced as the way to patch up this final difficulty (as well as all the ones encountered previously) and, simultaneously, to invoke the “common knowledge” principle for beliefs. Comparison with Lemma 1 is useful: We require not only convergence of the probability distribution on endpoints to the distribution given by π (as in Lemma 1), but convergence to the strategy π everywhere, including those parts of the tree that could be reached only if an initial defection occurred. This is equivalent to lexicographic consistency where we add restrictions on how players devise their secondary hypotheses $\pi(2), \pi(3), \dots$; $\pi(2)$ should be a “minimal change” from π , and so forth. To state precisely the equivalent lexicographic characterization of the consistency criterion is quite involved, and we have relegated it to the first part of the Appendix. The interested reader may wish at this point to read that part of the Appendix, in order to understand how our consistency criterion formalizes two ideas: no correlation in defections from π ; common knowledge in secondary hypotheses and in the formulation of beliefs.

Upon studying the first part of the Appendix, the reader may well conclude that we have required too much consistency in beliefs off the equilibrium path. Certainly, weaker consistency criteria can be posed that make sense in particular contexts. (But beware: By dropping (5.4), subgame perfection is lost.) We shall proceed here to develop the properties of sequential equilibrium as defined above; however, we do so with some doubts of our own concerning what “ought” to be the definition of a consistent assessment that, with sequential rationality, will give the “proper” definition of a sequential equilibrium.

6. PROPERTIES OF SEQUENTIAL EQUILIBRIA

We begin by establishing that standard properties hold for sequential equilibria.

PROPOSITION 1: *For every extensive game, there exists at least one sequential equilibrium.*

PROPOSITION 2: *Fixing an extensive form, the correspondence from pairs (ρ, u) of initial assessments and payoffs to the set of sequential equilibria for the game so defined is upper hemi-continuous.*

Proposition 1 is an easy corollary of Proposition 5 and Theorem 5 of Selten [15]. The proof of Proposition 2 follows the usual lines and is left to the reader.

The next result is no more than a marshalling of definitions.

PROPOSITION 3: *If (μ, π) is a sequential equilibrium, then π is a subgame perfect Nash equilibrium.*

Recall that subgame perfection is limited to proper subforms because for general subforms \hat{T} the specification of (ρ, u) and a strategy π does not necessar-

ily yield an initial assessment $\hat{\rho}$ on the set of initial nodes \hat{W} of \hat{T} . If the given strategy π is such that $P^\pi(\hat{W}) > 0$, then the appropriate choice for $\hat{\rho}$ is manifestly $\hat{\rho}(\hat{w}) = P^\pi(\hat{w})/P^\pi(\hat{W})$; in this case it is easy to show that if π is a Nash equilibrium then π restricted to \hat{T} is a Nash equilibrium on the subgame given by \hat{T} , $\hat{\rho}$, and u . Of course, in general $\hat{\rho}$ will not be strictly positive.

In seeking to extend this construction to general subforms we make the following definition.

DEFINITION: The strategy π is *extended subgame perfect* if for every subform \hat{T} there exists a nonnegative probability measure $\hat{\rho}$ on \hat{W} such that together with ρ on \bar{W} , (a) π is a Nash equilibrium for the game defined from \hat{T} , $\hat{\rho}$ and u ; and (b) if $\bar{T} \subseteq \hat{T}$ are two subforms such that $P^{\pi, \hat{\rho}}(\bar{T} | \hat{T}) > 0$, then $\bar{\rho}$ is defined from π and $\hat{\rho}$ by Bayes' formula.

(We have not formally defined $P^{\pi, \hat{\rho}}(\cdot | \hat{T})$, but its definition should be apparent.) Part (b) of this definition can be paraphrased: The initializing measures $\hat{\rho}$ are related to each other (and to ρ) in the natural fashion, given π .

PROPOSITION 4: *If (μ, π) is a sequential equilibrium, then (μ, π) is extended subgame perfect.*

The proof is left to the reader. The basic idea is quite simple: Suppose that (μ, π) is a sequential equilibrium. Let $\{(\mu_n, \pi_n)\}$ be some sequence from Ψ^0 with limit (μ, π) . Then for a subform \hat{T} with initial nodes \hat{W} , define $\hat{\rho}(\hat{w}) = \lim_n P^{\pi_n}(\hat{w})/P^{\pi_n}(\hat{W})$ for $\hat{w} \in \hat{W}$, looking along a subsequence if necessary. It is straightforward to verify that π is a Nash equilibrium for the game so defined.

The reader may wonder whether something of a converse of Proposition 4 is possible. Namely, if π is extended subgame perfect, is it then part of a sequential equilibrium? The answer is no, and an example is given in Figure 8. The information sets in this game are constructed so that the only subform is the original form. Thus every Nash equilibrium is extended subgame perfect. The strategies indicated form a Nash equilibrium, but they could never be part of a sequential equilibrium: No matter what beliefs player 2 has in his information set, B is a better action than A . (And if 2 plays B , then 1 and 3 prefer b and b' —the unique sequential equilibrium has the strategies depicted completely reversed.)

We turn next to a study of the topological structure of the set of sequential equilibria for generic games. The analysis follows Debreu's [1] study of the Walrasian correspondence.

Fix an arbitrary extensive form and initial assessment ρ . To specify an extensive game, it remains to specify a payoff vector $u \in R^{I \times Z}$. We say that a statement is true generically or for generic u if the closure of the subset in $R^{I \times Z}$ for which it is false has Lebesgue measure zero.

Let $\Phi(u)$ denote the set of sequential equilibria for the game with payoffs u . In general, the sets $\Phi(u)$ can be quite complicated. Two examples will illustrate this.

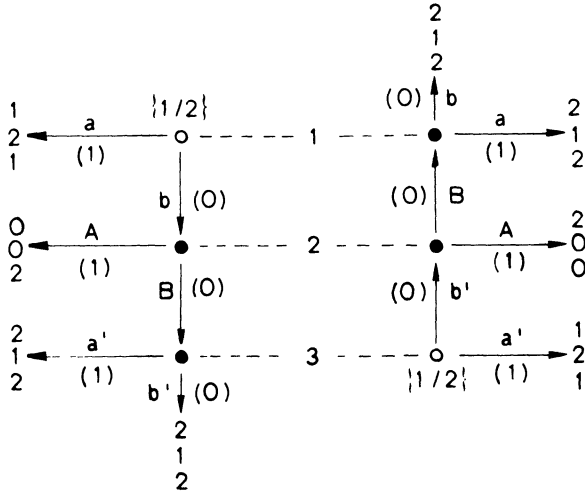


FIGURE 8.

Consider the game depicted in Figure 9. This game has two types of sequential equilibria. The first has player 1 playing L with probability one and 2 playing l with probability one. The second has 1 playing A with probability one and 2 playing r with probability $1/3$ or more. The latter type requires 2 to have beliefs assigning x a probability not exceeding $1/2$. If we project $\Phi(u)$ into the space of pairs $(\mu(x), \pi(l))$, we have the picture given in Figure 10. Note well the isolated point $(1, 1)$.

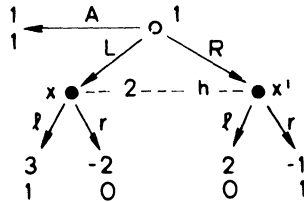


FIGURE 9.

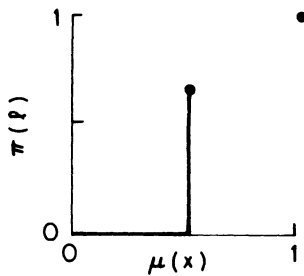


FIGURE 10.

The second example concerns the game depicted in Figure 11. It is easiest to think of this game as follows. Player 1 either chooses *A* or chooses one of the following two bimatrix games for 2 and 3 to play:

		3	
		<i>l'</i>	<i>r'</i>
2	<i>l</i>	1, 0	0, 1
	<i>r</i>	0, 1	1, 0
		(<i>L</i>)	

		3	
		<i>l'</i>	<i>r'</i>
2	<i>l</i>	1, 1	0, 0
	<i>r</i>	0, 1	0, 0
		(<i>R</i>)	

Neither 2 nor 3 knows which bimatrix game is selected. Again there are two types of equilibria. In the first, 1 moves *R* with probability one, and 2 and 3 respond with *l* and *l'*. In the second type, 1 moves *A* with probability one. This gives 2 and 3 “freedom” in their beliefs concerning whether 1 chose *L* or *R* if information sets *h* and *h'* are reached. (Because we are looking for beliefs that are consistent, these beliefs must coincide.) To have a sequential equilibrium, these beliefs must satisfy $\mu(x) \in [1/2, 1/1.1]$, with resulting equilibrium strategies

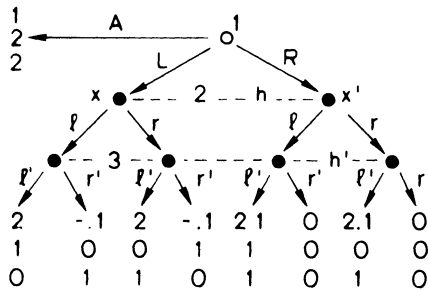


FIGURE 11.

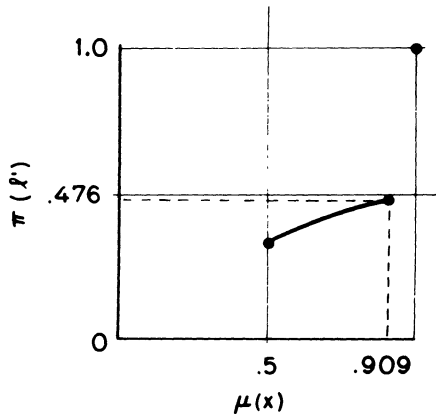


FIGURE 12.

$\pi(l) = 1/(2\mu(x))$ and $\pi(l') = \mu(x)/(1 + \mu(x))$. Projecting these equilibria into the space of pairs $(\mu(x), \pi(l'))$ we get the set shown in Figure 12. Note that the isolated point corresponds to the first type of equilibrium, and the curved segment to the second.

In both examples, small perturbations in u do not affect the basic shape of $\Phi(u)$ —the same shape holds for all vectors u in a neighborhood of the payoffs depicted.

These examples illustrate the structure of $\phi(u)$ for generic u : $\phi(u)$ is typically the union of manifolds of various dimensions. The dimensions of these manifolds are related to the number of “degrees of freedom” that are available in specifying beliefs or strategies off the equilibrium path. The easiest case to understand is illustrated by the horizontal segment in Figure 10 and the curved segment in Figure 12. In $\Phi(u)$ (as well as in the projections illustrated) these segments are one-dimensional manifolds derived from the one degree of freedom available to specify the non-Bayesian beliefs $\mu | h$. In Figure 12 we see that this specification may have some effect on the corresponding equilibrium strategies. The vertical segment in Figure 10 has a more subtle explanation: The “degree of freedom” that exists in defining μ on h is lost because μ on h is set so that player 2 is indifferent between l and r , but this degree of freedom is regained in the choice of π on h because player 2 is indifferent.

Several definitions are required to be precise. A *basis* for the extensive form is an index set b consisting of decision nodes $x \in X$ and actions $a \in A$. Define $\Psi_b = \{(\mu, \pi) \in \Psi : \mu(x) > 0 \text{ if and only if } x \in b, \text{ and } \pi(a) > 0 \text{ if and only if } a \in b\}$. The set of bases for which Ψ_b is nonempty is denoted by B and is called the set of *consistent* bases. Note that B is finite and that the Ψ_b partition Ψ .

LEMMA 2: For $b \in B$, Ψ_b is a manifold.

This lemma is proved in the Appendix. Also in the Appendix is a characterization of which bases b are consistent, and a representation of the manifold Ψ_b for consistent b .

Partition $\Phi(u)$ as follows. Define $\Phi_b(u) = \Phi(u) \cap \Psi_b$. Further partition each $\Phi_b(u)$ into two parts: $\Phi_b^S(u)$ —the set of *strict equilibria* in Ψ_b , wherein any action that does as well as an action taken with positive probability is itself taken with positive probability; $\Phi_b^W(u)$ —the remainder of $\Phi_b(u)$, consisting of the *weak equilibria* wherein some unused action does as well as the actions having positive probability.

THEOREM 1: For generic u , for each $b \in B$ the set $\Phi_b^S(u)$ is either empty or is a manifold of dimension $n(b) = \dim(\Psi_b) - \#(b \cap A) + \#H$, and the set $\Phi_b^W(u)$ is precisely the Ψ_b -relative frontier of $\Phi_b^S(u)$.

TABLE II
 $\Phi_b(u)$ FOR THE EXAMPLE OF FIGURE 9

basis b	$\dim(\Psi_b)$	$n(b)$	$\Phi_b^S(u)$	$\Phi_b^W(u)$
$A; x'; r$	0	0	the point at the origin in Figure 9	\emptyset
$A; x, x'; r$	1	1	the horizontal segment not including either endpoint	the right endpoint of the horizontal segment
$A; x, x'; r, l$	2	1	the vertical segment not including either endpoint	the top endpoint of the vertical segment
$L; x; l$	0	0	the isolated point at the upper right	\emptyset

The proof is given in the Appendix. An intuitive explanation of the dimension of $\Phi_b^S(u)$ can be given: Fixing a basis b and an information set h , suppose that m is the cardinality of $A(h) \cap b$. Then the equilibrium conditions for h consist of $m - 1$ equality constraints specifying that the expected utilities to player $i(h)$ of the m actions he is mixing among are identical, and some further inequality constraints. Each equality constraint lowers the dimension of the manifold of equilibria by one, so when we sum over all h we have exactly $n(b)$ remaining dimensions of Ψ_b .

To see how this theorem works, consider again the game depicted in Figure 9. The bases $b \in B$ such that $\Phi_b^S(u)$ is nonempty are tabulated in the left column of Table II. Also tabulated are $\dim(\Psi_b)$ and $n(b)$. In the next column, $\Phi_b^S(u)$ is described, relying on Figure 9. And in the right column, the associated $\Phi_b^W(u)$ is described.

This theorem points out that typically there is an infinite number of sequential equilibria. This may be misleading, however. The infinite number ensues from free variations in allowed behavior and beliefs *off* the equilibrium path. To be more precise, for a fixed extensive form and initial assessment ρ , define

$$\Delta_u = \{P^\pi(\cdot) \in \Delta : \text{There exists some } \mu \text{ such that } (\mu, \pi) \in \Phi(u)\}.$$

In a sense, Δ_u is the projection of $\Phi(u)$ onto Δ . Note that a measure $\delta \in \Delta_u$ represents an equilibrium prior probability assessment on the terminal node at which the game will end.

THEOREM 2: *For generic u , the set Δ_u is finite.*

The proof is given in the Appendix. Also in the Appendix are further characterizations of the topological properties of the correspondences $\Phi(u)$ and Δ_u . We add two comments here. First, Theorem 2 remains true if we look at all Nash equilibria and not only at sequential equilibria. Second, for generic normal form games the set Δ_u has odd cardinality. In contrast, this is not the case for generic extensive games; cf. the two examples of this section.

7. PERFECT EQUILIBRIA

In this section we establish the relation between sequential and perfect equilibria. The reader familiar with Selten's [15] exposition of the merits of perfect equilibria will recognize that our justification for the criterion of sequential rationality substantially overlaps his. Indeed, our analysis has relied heavily on the creative insights of Selten's original work, and on the related contributions of Harsanyi [4]. We have nevertheless proposed a definition that is formulated quite differently, and that yields somewhat different properties. Our approach involves directly the criterion of sequential rationality, whereas Selten employs an indirect construction to obtain a slightly stronger criterion. Our motivation for simplifying Selten's construction is principally a matter of analytic ease. Sequential equilibria are generally much easier to work with. Below we shall recall Selten's definition, contrast it with our own, and then show the sense in which the two are virtually equivalent.

To facilitate comparisons we formulate Selten's definition in terms of assessments. (Also, we shall follow Selten's second definition [15, Section 12] rather than his first.) Recall that the set Ψ of consistent assessments is the closure of the set Ψ^0 of strictly positive assessments. We say that a convergent sequence $\{(\mu_n, \pi_n)\}_{n=1,2,\dots}$ from Ψ^0 justifies the fully consistent assessment (μ, π) that is its limit. Then, an assessment $(\mu, \pi) \in \Psi$ is a *perfect equilibrium* if it is justified by a sequence $\{(\mu_n, \pi_n)\}$ for which, for each player i and each index n , π^i is an optimal response for player i to the other players' strategies $(\pi_n^j)_{j \neq i}$.

One can interpret this definition as the composition of two criteria. The first, implied by the above definition, is that each player i 's strategy π^i is an optimal response to the assessment (μ, π) and that, moreover, (μ, π) is a fully consistent sequential equilibrium. The second requires that for some sequence that justifies (μ, π) , each player's strategy is a robust best response. One can interpret the difference between a strategy π_n in the sequence and the limit strategy π as the manifestation of particular small probabilities that the other players will err or "tremble", and that each player is using a response that is optimal in the event of such errors. It is the second of these two criteria that we forego in the definition of a sequential equilibrium. This makes apparent the following.

PROPOSITION 5: *Every perfect equilibrium is a sequential equilibrium.*

The converse to this is false. In the game depicted in Figure 13, player 1 moving L and 2 moving r is a sequential equilibrium that is not perfect. Moreover, one can find more complicated games with weak sequential equilibria that are not perfect for all payoffs u in a neighborhood of some u^* . However, the following partial converse to Proposition 5 is true.

THEOREM 3: *For any fixed extensive form and initial assessment ρ , for generic payoffs u every strict sequential equilibrium is perfect. Also, for generic payoffs u the projection from perfect equilibria to Δ coincides with Δ_u .*

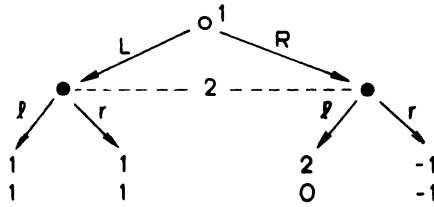


FIGURE 13.

In other words, for generic payoffs, only weak sequential equilibria present a problem, and then only off the equilibrium path. From Theorem 1, we know that (generically) “almost every” sequential equilibrium is strict, so we can further paraphrase Theorem 3 as follows: For almost every game, almost every sequential equilibrium is perfect. The proof is in the Appendix. It is a corollary of the construction in the Appendix, moreover, that the sets of sequential and perfect equilibria fail to coincide only at points u where the perfect equilibrium correspondence fails to be upper hemi-continuous.

By modifying Selten’s definition of perfection, we can obtain exact coincidence of sequential and “perfect” equilibria. To motivate this, we embellish the motivation for perfect equilibria: Imagine that some player i is predicting the behavior of another player j . To do this requires that i consider j ’s predictions of the behavior of everyone else. That is, the Nash criterion requires that j ’s behavior π^j be a best response to the particular prediction $(\pi^k)_{k \neq j}$ of others’ behavior. Selten has added that j ’s behavior π^j must also be a best response to some sequence of strictly positive strategies $(\pi_n^k)_{k \neq j}$ that approach π . The interpretation is that if some information set $h \in \iota^{-1}(j)$ is reached that has zero prior probability under π , then j will interpret this event to result from some (vanishingly) small chance of a sequence of “errors” and proceed to optimize accordingly. That is, player i must suppose that j ’s actions are a robust best response. This all supposes that i knows j ’s payoffs. We can slightly relax Selten’s criterion by allowing some (vanishingly) small uncertainty on the part of i about j ’s payoffs; then j ’s strategy need only be a best response to the perturbed strategies for some payoffs for j that are “close” to u . Formally we have the following definition.

DEFINITION: An assessment $(\mu, \pi) \in \Psi$ is said to be a *weak perfect equilibrium* for payoffs u if there exists a sequence $\{(\mu_n, \pi_n, u_n)\} \subseteq \Psi^0 \times R^{1 \times Z}$ that has the limit (μ, π, u) and that satisfies, for each n and player j , π^j is a best response to $(\pi_n^k)_{k \neq j}$ if j ’s payoffs are u_n^j .

PROPOSITION 6: *For any extensive game, the sets of weak perfect and sequential equilibria coincide.*

The proof is quite straightforward, so we only sketch the main idea and leave details to the reader. Take any sequence $\{(\mu_n, \pi_n)\}$ from Ψ^0 with limit (μ, π) . Fix

an n . Then working backwards through the game tree, it is possible to perturb slightly each u^j so that π^j is a best response to $(\pi_n^k)_{k \neq j}$ for the perturbed payoffs. (See Part 2 of the Appendix for further details on this backwards recursion.)

While we are concerned in this paper with extensive games, our three theorems combine to give an interesting corollary for normal form games. In a normal form game, every Nash equilibrium is sequential. Moreover, knowing P^π is equivalent to knowing π . Thus we conclude that for generic normal form games, there is a finite number of Nash equilibria, every one of which is perfect.

8. RESTRICTIONS ON BELIEFS

Besides being easier to apply than perfectness, the concept of sequential equilibrium possesses a second advantage. This is that the formulation in terms of players' beliefs gives the analyst a tool for choosing among sequential equilibria. In some cases one can predict that one equilibrium among several will prevail because only the beliefs that sustain the one are intuitively plausible. An example illustrates this.

Consider the game depicted in Figure 14, due to Kohlberg [5]. One sequential equilibrium for this game has 1 playing A and 2 playing r . This is supported by beliefs on the part of 2 that $\mu(x) \leq 1/2$. A second equilibrium has 1 playing L and 2 playing l . In this equilibrium, 2's beliefs are determined by 1's strategy. We contend that the second equilibrium is more intuitively plausible than the first, because 2's beliefs in the first are implausible. Player 2 "ought" not to conclude, upon reaching h , that 1 chose R , because R is dominated by A for 1. The only beliefs by 2 that make sense assign $\mu(x) = 1$ (or, at least, $> 1/2$), which leads 2 to prefer l . Player 1, realizing this, prefers L .

Now consider the following modification of this game. Suppose that by choosing A , player 1 causes the following simultaneous-move, bimatrix game to be played:

		2		
		2, 1	0, 2	- 10, 4
1	L	0, 4	2, 2	- 8, 1

Player 2 will know that 1 chose A and that they are therefore playing this game. In this bimatrix game, there is a unique Nash equilibrium where 1 receives the

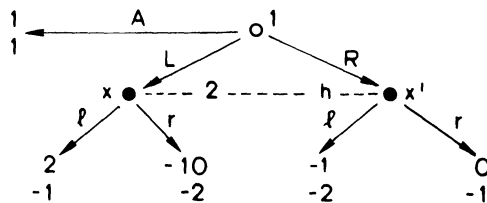


FIGURE 14.

expected payoff 1. Again we have a sequential equilibrium in the whole game where 1 chooses A because 2, upon reaching h , can assess $\mu(x) < 1/2$ and therefore play r . And again we argue that this is implausible because 2's beliefs are implausible. It is no longer the case that playing A guarantees 1 more than he can ever get by playing R . But playing A yields an expected equilibrium payoff of 1, which is more than 1 gets from R , so 2 "ought not" to attach too much weight to the possibility that 1 played R given that 1 has not played A . Again we are led to the equilibrium L and l as more plausible.

One can also interpret Myerson's [10] concept of a *proper* equilibrium as a restriction on beliefs. Myerson argues (essentially) that assessments off the equilibrium path should be such that the preponderance of weight goes to the "least costly" mistakes. A third example of this approach is used in Kreps and Wilson [6, Section 3]; there assessments off the equilibrium path are developed by asking which of two players is *more* likely to defect from a given equilibrium.

The point that we wish to make is that the basic description, and justifying feature, of sequential equilibria, namely probabilistic reassessments of beliefs off the equilibrium path, provide an apparatus for comparing sequential equilibria. Some sequential equilibria are supported by beliefs that the analyst can reject because they are supported by beliefs that are implausible. We will not propose any formal criteria for "plausible beliefs" here. In certain cases, such as Myerson's concept of properness, some formalization is possible. In other cases it is not clear that any formal criteria can be devised—it may be that arguments must be tailored to the particular game. But whether this is done formally or informally, we believe that this perspective can occasionally be employed to advantage.

9. CONCLUDING REMARKS

The criterion of sequential rationality is familiar in the analysis of single-person decision problems. It justifies the standard "roll-back" procedure for constructing a sequentially optimal strategy in a problem described by a decision tree. Extensive games are multi-person decision problems. Though complicated by the interactions among players, they are not different in substance. An extensive game with perfect recall factors into decision problems (and associated decision trees) for each player. The complicating feature is that the events in one player's decision tree may correspond to actions of others. A player's probability assessments must, therefore, depend upon anticipations of others' strategies. An equilibrium in Nash's sense supposes that strategies are "common knowledge" among the players. Consequently, strategies and beliefs are intertwined in complicated ways by the information structure of the game. To satisfy sequential rationality one must use them in concert to verify an equilibrium.

We have formulated the concept of sequential equilibrium to emphasize this view. Compared to the weaker criterion of Nash equilibrium, the salient difference is the key role of players' beliefs off the equilibrium path in determining optimal strategies subsequent to unanticipated events. Compared to the slightly

stronger criterion of trembling-hand perfect equilibrium, the difference is mainly one of tractability. The simpler mathematical properties of the sequential equilibrium correspondence reflect the simpler formulation of the basic criterion.

We envision the criterion of sequential equilibrium as a step towards unifying the formulations and methodologies of game theory and statistical decision theory. Among the prospective applications are analyses of games modeling competitive strategies in dynamic and uncertain environments. The criterion of sequential rationality subjects the (usually) many Nash equilibria to a stronger test. As Selten has noted, it is particularly effective in eliminating equilibria sustained by "threats" that would be nonoptimal to carry out.

The idea that an assessment, namely a beliefs-strategy pair, is the relevant descriptor of an equilibrium is, we contend, a central principle. Analyses that ignore the role of beliefs, such as analysis based on normal-form representation, inherently ignore the role of anticipated actions off the equilibrium path in sustaining the equilibrium—essentially, such analyses allow almost arbitrary behavior off the path. This lacuna often weakens the normative implications of the analysis, and in the extreme yields Nash equilibria that are patently implausible. This is not to say that the sequential equilibrium concept eliminates all or even most implausible equilibria. But it does eliminate some, and, perhaps more importantly, it gives a language for discussing why one equilibrium or another is implausible.

We close by noting that the ideas here are not original to us. We have already cited the seminal work of Harsanyi and Selten, but it would be remiss not to do so again. Selten's work on perfection goes at precisely the problem we began with. Our concept of sequential equilibrium allows a somewhat freer interpretation of the nature of beliefs off the equilibrium path, in that Selten motivates them entirely by "trembling-hands." But mathematically the two are equivalent (for fully consistent assessments). (The other distinction, that we require optimality only "at the limit," while Selten requires optimality approaching the limit, is what is significant in terms of tractability and mathematical properties.) And our freer interpretation springs directly from Harsanyi's work on games of incomplete information—we simply apply those ideas to players' decision tree problems off the equilibrium path. Also, recent papers by Fudenberg and Tirole [2], Milgrom and Roberts [8], and Rubinstein [12] contain the basic idea of a sequential equilibrium without being quite so formal about it as we have been.

Stanford University

Manuscript received February, 1981; revision received June, 1981.

APPENDIX

In this Appendix we prove Lemma 2 and Theorems 1, 2, and 3. We begin by considering the structure of the sets Ψ_b , developing characterizations of consistent bases, and an explicit representation of Ψ_b . Lemma 2 is then proven, and an extension is given that will be used in the proof of Theorem 3. In the second part of the Appendix, we present a general construction that is used

repeatedly in the proofs of the theorems. Then in the third part of the Appendix, we give proofs of the three theorems. The basic mathematical tool of our analysis is Sard's Theorem, used in the fashion initiated by Debreu [1].

A.1. THE STRUCTURE OF Ψ_b

To simplify formulae, we assume throughout this section that $\rho \equiv 1/\#W$. This is for notational convenience only—none of the analysis to follow depends on this at all.

Begin with a definition. A *labelling* for the extensive form is a function K taking A into the nonnegative integers. For a given labelling K , there is an associated function J_K on X defined by

$$J_K(x) = \sum_{l=0}^{l(x)-1} K(\alpha(p_l(x))).$$

If $x \in W$, then set $J_K(x) = 0$. That is, K labels the branches of the tree with nonnegative integers (in a way that respects the informational constraints of the game) and J_K gives for each node x the sum of the labels on branches from the beginning of the tree to x .

The labelling K is said to be a *b labelling* if:

- (A.1) (a) for every h , there is some $a \in A(h)$ with $K(a) = 0$;
- (b) $a \in b$ if and only if $K(a) = 0$;
- (c) $x \in b$ if and only if x minimizes $J_K(\cdot)$ on $H(x)$.

LEMMA A1: *The basis b is consistent (Ψ_b is nonempty) if and only if a b labelling exists.*

PROOF: Suppose that a b labelling K exists. Fix any strategy $\pi \in \Pi^0$ and define strategies $\pi_n \in \Pi^0$ by

$$\pi_n(a) = c(n, H(a))\pi(a)(1/n)^{K(a)},$$

where $c(n, H(a))$ is defined as the appropriate normalizing constant. Letting μ_n be the beliefs consistent with π_n , it is obvious that the sequence $\{(\mu_n, \pi_n)\}$ converges to some assessment (μ, π) that belongs to the basis b . Thus Ψ_b is nonempty.

Now suppose that b is a consistent basis. Since Ψ_b is nonempty, there exists a sequence $\{(\mu_n, \pi_n)\} \subseteq \Psi^0$ with the limit (μ, π) belonging to Ψ_b . Let M denote the finite set of all first degree, single term multinomials with coefficient one in the symbols $a \in A$. For $m \in M$, let m_n represent m evaluated with $a = \pi_n(a)$. Without loss of generality, we can assume that for every pair m and m' from M , the sequence m_n/m'_n converges either to zero, to infinity, or to some strictly positive number. (This is wlog because we can look along a subsequence of $\{(\mu_n, \pi_n)\}$ for which it is true.) Define $m < m'$ if $\lim_n m_n/m'_n = \infty$; then $<$ is an asymmetric and negatively transitive binary relation on M . Since M is finite there exists an integer valued function J on M with $m < m'$ if and only if $J(m) < J(m')$. We can pick J so that $J(m) = 0$ for the $<$ -least m —then $J(m) \geq 0$ for all m . For each $x \in X$ there is an associated $m^x \in M$, namely

$$m^x = \prod_{l=0}^{l(x)-1} \alpha(p_l(x)).$$

(For $x \in W$, $m^x = 1$.) Now for each a pick an arbitrary $x \in H(a)$ such that $J(m^x)$ is minimal over $x \in H(a)$ and define

$$K(a) = J(m^x \cdot a) - J(m^x).$$

We leave to the reader the relatively easy tasks of proving that $K(a)$ is well-defined (i.e., the choice of a $J(m^x)$ -minimal $x \in H(a)$ is irrelevant) and that K so defined is a b labelling (with, of course, $J_K(x) = J(m^x)$). Q.E.D.

We can now give the representation of Ψ_b . Let Ξ_b be the set of functions $\xi : A \rightarrow (0, \infty)$ such that for each $h \in H$,

$$\sum_{a \in b \cap A(h)} \xi(a) = 1.$$

That is, each ξ is a strategy that belongs to the basis b together with an assignment of positive numbers to all other branches. Next define mappings π^b and μ^b from Ξ_b to Π and the space of beliefs, respectively, by

$$\pi^b(\xi)(a) = \begin{cases} 0 & \text{if } a \notin b, \text{ and} \\ \xi(a) & \text{if } a \in b, \end{cases}$$

and

$$\mu^b(\xi)(x) = \begin{cases} 0 & \text{if } x \notin b, \\ m^x(\xi) / \sum_{x' \in b \cap H(x)} m^{x'}(\xi) & \text{if } x \in b, \end{cases}$$

where $m^x(\xi)$ is the multinomial m^x evaluated with $a = \xi(a)$.

LEMMA A2: For each consistent basis b , Ψ_b is the image of Ξ_b under the mapping (μ^b, π^b) .

PROOF: Fix a consistent basis b . We first show that for $\xi \in \Xi_b$, $(\mu^b(\xi), \pi^b(\xi)) \in \Psi_b$. To do this, let K be any b labelling (one exists by Lemma A1) and, for $n = 1, 2, \dots$, define $\pi_n(\xi) \in \Pi^0$ by

$$\pi_n(\xi)(a) = c(n, H(a))\pi(a)(1/n)^{K(a)},$$

where $c(n, H(a))$ is the appropriate normalizing constant. If $\mu_n(\xi)$ are the beliefs consistent with $\pi_n(\xi)$, it is easy to see that

$$\lim_{n \rightarrow \infty} (\mu_n(\xi), \pi_n(\xi)) = (\mu^b(\xi), \pi^b(\xi)),$$

so $(\mu^b(\xi), \pi^b(\xi)) \in \Psi$. (Note that each $c(n, H(a))$ goes to one.) By definition, $(\mu^b(\xi), \pi^b(\xi))$ has the basis b and hence is in Ψ_b .

Conversely, we must show that for every $(\mu, \pi) \in \Psi_b$ there exists $\xi \in \Xi_b$ with $(\mu, \pi) = (\mu^b(\xi), \pi^b(\xi))$. Fix $(\mu, \pi) \in \Psi_b$, and let $\{(\mu_n, \pi_n)\} \subseteq \Psi^0$ be any fixed sequence that justifies (has limit) (μ, π) . Let \mathcal{Q} be the set of all algebraic expressions of the form

$$q = \prod_{a \in A} a^{q(a)}$$

for rational numbers $q(a)$. Note that \mathcal{M} is the subset of \mathcal{Q} for which each $q(a)$ is either zero or one. Note that \mathcal{Q} has a countable number of elements. Without loss of generality, then, we can assume that for each pair $q, q' \in \mathcal{Q}$ the limit $\lim_n q_n/q'_n$ exists (allowing ∞ as a limit), where q_n is the expression q evaluated with $a = \pi_n(a)$. This is wlog because we can use the standard diagonalization procedure to find a subsequence along which it is true. Define binary relations $<$ on \mathcal{Q} as before, and let \doteq represent the associated equivalence relation.

Now we define $\xi(a)$. For $a \in b$, set $\xi(a) = \pi(a)$. We shall define the remaining $\xi(a)$ one at a time, taking care so that at each step along the way the following statement is true:

If q and q' are from \mathcal{Q} with $q(a) = q'(a) = 0$ for any a whose $\xi(a)$ value has not yet been assigned, and if $q \doteq q'$, then $\lim_n q_n/q'_n = q(\xi)/q'(\xi)$.

Here $q(\xi)$ means q evaluated with $a = \xi(a)$. To verify that this is possible, note first that for the initial assignment given above, the statement is true for $a \in b$ —this follows from $\lim_n \pi_n = \pi$ and $\pi(a) > 0$ for $a \in b$. Now suppose that we have assigned $\xi(a)$ for some subset of A so that the statement holds, and we try to assign $\xi(a^*)$ for some unassigned a^* so that it holds for the augmented subset.

If q' and q'' are two algebraic expressions in the already-assigned a and in a^* such that $q' \doteq q''$, then we can formally solve the “equation” $q' = q''$ for a^* to get $a^* = q$ for some q in the

already-assigned a . Moreover, it is easy to see that for this q , $a^* \doteq q$. And if we ensure that $\lim_n \pi_n(a^*)/q_n = \xi(a^*)/q(\xi)$ for this q , then we will have $\lim_n q_n/q_n = q'(\xi)/q''(\xi)$ for the original q' and q'' . So in checking that an assignment of $\xi(a^*)$ can be made that preserves the statement, we need only check the cases $a^* \doteq q$ for q an expression in the already-assigned a .

Now if *no* q in the already-assigned a satisfies $a^* \doteq q$, $\xi(a^*)$ can be arbitrarily assigned. While if $a^* \doteq q$ and $a^* \doteq \hat{q}$, we have $q \doteq \hat{q}$, and the "induction hypothesis" implies that $\lim_n q_n/\hat{q}_n = q(\xi)/\hat{q}(\xi)$, which in turn implies that the necessary assignment of $\xi(a^*)$,

$$\xi(a^*) = \left[\lim_n \pi_n(a^*)/q_n \right] \cdot q(\xi),$$

does not depend on which $q \doteq a^*$ is selected.

All this implies that we can find an assignment of $\xi(a)$ (with $\xi(a) = \pi(a)$ for $a \in b$) so that the statement holds for *all* $q, q' \in Q$. In particular, the statement holds for each pair m^x and $m^{x'}$ where x and x' are both in b and are in the same information set. For such x and x' we know that $\mu(x)$ and $\mu(x')$ are both greater than zero, and thus that $m^x \doteq m^{x'}$. Therefore, $\lim_n (m^x)_n / (m^{x'})_n = m^x(\xi) / m^{x'}(\xi)$ which immediately implies that $\mu = \mu^b(\xi)$. That $\pi = \pi^b(\xi)$ is apparent, and we are done. Q.E.D.

PROOF OF LEMMA 2: We now make a convenient change of variables. Fixing a consistent basis b , let Z_b be the subset of $\zeta \in R^A$ satisfying

$$\sum_{a \in b \cap A(h)} \exp(\zeta(a)) = 1 \quad \text{for every } h \in H.$$

Let $e : Z_b \rightarrow \Xi_b$ be the map $e(\zeta)(a) = \exp(\zeta(a))$. Clearly, e is an isomorphism. Now redefine μ^b and π^b above so that their domain is Z_b . That is, set $\pi^b(\zeta)(a) = \exp(\zeta(a))$ for $a \in b$, and so on. By Lemma A2, Ψ_b is the image of Z_b under the (redefined) map (μ^b, π^b) .

Define

$$\Theta_b = \left\{ \theta \in R^A : \theta(a) = 0 \text{ if } a \in b; \text{ and for each } h \text{ if } x, x' \in b \cap h \text{ then } \sum_{l=0}^{l(x)-1} \theta(\alpha(p_l(x))) = \sum_{l=0}^{l(x')-1} \theta(\alpha(p_l(x'))) \right\}.$$

Note that Θ^b is a subspace of R^A and that $\Theta^b + Z^b = Z^b$. Moreover, we assert that the null space of the Jacobian of (μ^b, π^b) is Θ^b at every point $\zeta \in Z^b$. And moreover, $(\mu^b, \pi^b)(\zeta) = (\mu^b, \pi^b)(\zeta')$ if and only if $\zeta - \zeta' \in \Theta^b$.

Some temporary notation will be useful. For $\omega \in R^A$ and $x \in X$, let $\Sigma(x, \omega)$ denote $\sum_{l=0}^{l(x)-1} \omega(\alpha(p_l(x)))$. Then for $\zeta \in Z_b$ and $x \in b, h = H(x)$, we have

$$\mu^b(\zeta)(x) = \exp(\Sigma(x, \zeta)) / \sum_{x' \in b \cap h} \exp(\Sigma(x', \zeta)).$$

So for arbitrary $\zeta \in Z_b$ and $\theta \in \Theta_b$ we have

$$\begin{aligned} \mu^b(\zeta + \theta)(x) &= \exp(\Sigma(x, \zeta + \theta)) / \sum_{x' \in b \cap h} \exp(\Sigma(x', \zeta + \theta)) \\ &= \exp(\Sigma(x, \zeta) + \Sigma(x, \theta)) / \sum_{x' \in b \cap h} \exp(\Sigma(x', \zeta) + \Sigma(x', \theta)) \\ &= \frac{\exp(\Sigma(x, \zeta)) \cdot \exp(\Sigma(x, \theta))}{\sum_{x' \in b \cap h} \exp(\Sigma(x', \zeta)) \exp(\Sigma(x', \theta))} = \frac{\exp(\Sigma(x, \zeta))}{\sum_{x' \in b \cap h} \exp(\Sigma(x', \zeta))} \\ &= \mu^b(\zeta)(x). \end{aligned}$$

Of course, $\pi^b(\zeta + \theta) = \pi^b(\zeta)$ for all $\zeta \in Z_b$ and $\theta \in \Theta_b$, since $\theta(a) = 0$ for $a \in b$. Thus $(\mu^b, \pi^b)(\zeta)$

$= (\mu^b, \pi^b)(\zeta')$ if $\zeta - \zeta' \in \Theta_b$. This in turn implies that Θ_b is always contained in the null-space of the Jacobian of (μ^b, π^b) .

Conversely, suppose that $\zeta - \zeta' = \delta \notin \Theta_b$. If $\delta(a) \neq 0$ for some $a \in b$, then $\pi^b(\zeta) \neq \pi^b(\zeta')$. And if for some $x \in b$, $\Sigma(x, \delta) \neq \Sigma(x', \delta)$ for some $x' \in b \cap H(x)$, then letting x^* be that element of $b \cap H(x)$ with largest $\Sigma(\cdot, \delta)$, it is easy to see that $\mu^b(\cdot)(x^*)$ is strictly increasing along the line from ζ' to ζ . Certainly, then, $(\mu^b, \pi^b)(\zeta) = (\mu^b, \pi^b)(\zeta')$ only if $\zeta - \zeta' \in \Theta_b$. And certainly, the null-space of the Jacobian of (μ^b, π^b) contains no vector not in Θ_b . (This last statement requires that we show that $\mu^b(\cdot)(x^*)$ has nonzero derivative in the direction of ζ from ζ' . This is true, as the reader can verify.)

To complete the proof, note that the Jacobians of (μ^b, π^b) at two points ζ and ζ' such that $(\mu^b, \pi^b)(\zeta) = (\mu^b, \pi^b)(\zeta')$ are identical. This implies that $(\mu^b, \pi^b)(Z_b)$ is a manifold, with dimension equal to $\#A - \#H - \dim(\Theta_b)$. Q.E.D.

The proof of the lemma yields a small sharpening that is useful in the proof of Theorem 3. Let

$$\bar{Z}_b = Z_b \cap \text{perp}(\Theta_b).$$

Then it is clear from the proof of the lemma that (μ^b, π^b) is a diffeomorphism from \bar{Z}_b to Ψ_b .

LEMMA A3: For each consistent basis b we can find some manifold $\Psi_b^0 \subseteq \Psi^0$ such that $\Psi_b^0 \cup \Psi_b$ is a smooth manifold with boundary, the boundary being Ψ_b .

(For the definition of a smooth manifold with boundary, see Milnor [9].) We will only give the construction Ψ_b^0 here—verification of the lemma is straightforward. Fix a b labelling K and, for $\zeta \in \bar{Z}_b$ and $r \in (0, 1)$, define a strategy $\pi_{\zeta,r} \in \Pi^0$ by

$$\pi_{\zeta,r}(a) = c(r, H(a)) \exp(\zeta(r)) r^{K(a)}$$

for appropriate normalizing constants $c(r, H(a))$. Let $\mu_{\zeta,r}$ be the beliefs consistent with $\pi_{\zeta,r}$. Then for Ψ_b^0 , take all such $(\mu_{\zeta,r}, \pi_{\zeta,r})$ for $\zeta \in \bar{Z}_b$ and $r \in (0, 1)$.

A.2. A GENERAL CONSTRUCTION

In order to prove the three theorems, it will be convenient to work with the players' decision trees. It is easiest to think of the game tree as beginning with some initial node 0 that precedes all nodes. Then the decision tree $(T^i, <)$ of player i consists of all information sets $h \in H^i$, all actions $a \in A^i$, all terminal nodes $z \in Z$, and the node 0. Precedence is inherited in natural fashion from the game tree. (Note the reliance here on perfect recall.) The set of immediate successors to any node $y \in T^i \setminus Z$ in i 's decision tree will be denoted by $S^i(y)$. The immediate predecessor of $y \in T^i \setminus \{0\}$ will be denoted by $p^i(y)$. (See Wilson [6].)

For every $(\mu, \pi) \in \Psi$, we obtain a corresponding transition probability assessment on each player's decision tree as follows: Fixing the player i , the probability of transition from a node $h \in H^i$ to a node $a \in A^i$ where $a \in A(h) (= S^i(h))$ is simply $\pi^i(a)$. For a node $a \in A(h)$ where $h \in H^i$ and for $y \in S^i(a) (\subseteq H^i \cup Z)$, the transition probability from a to y is defined to be $P^{\mu, \pi(a)}(y | h)$, where $\pi(a)$ is the strategy π changed so that action a is taken with certainty in information set h . Transition probabilities from 0 to its immediate successors, each of which will be some $y \in H^i \cup Z$, are given by $P^\pi(y)$. Let v^i be the map from Ψ to the space of transition probabilities on player i 's decision tree, and let v denote the vector map $(v^i)_{i \in I}$. Clearly, v is a smooth map, consisting only of iterated multiplications and additions and having derivatives of all orders. We will write $v^i(\mu, \pi)(y)$ to denote the appropriate transition probability to a node $y \in T^i \setminus \{0\}$. Note that for each $y' \in T^i \setminus Z$, $\sum_{y \in S^i(y')} v^i(\mu, \pi)(y) = 1$ for all (μ, π) .

The following is simply a matter of marshalling definitions, and so is stated without proof.

LEMMA A3: For $(\mu, \pi) \in \Psi$, (μ, π) is a fully consistent sequential equilibrium for payoffs u if and only if for each i , when we iteratively compute

$$(A.2) \quad (a) \ v^i(a) = \sum_{y \in S^i(a)} v^i(\mu, \pi)(y) \cdot v^i(y), \quad \text{for } a \in A^i, \quad \text{and}$$

$$(b) \ v^i(h) = \max_{a \in A(h)} v^i(a) \quad \text{for } h \in H^i,$$

initializing with $v^i(z) = u^i(z)$, we find that a attains the maximum in (A.2)(b) if $\pi(a) > 0$. Moreover, a strict equilibrium is where a attains the maximum in (A.2)(b) if and only if $\pi(a) > 0$.

That is, an equilibrium in the game tree corresponds to “individually rational behavior” in the decision tree, using the standard roll-back procedure from decision analysis.

For each player i , let $g^i: T^i \setminus Z \rightarrow T^i \setminus \{0\}$ be any arbitrarily selected function that specifies for each nonterminal node y in i 's decision tree some immediate successor $g^i(y)$ of this node. Define

$$D^i = \{d \in R^T : d(g^i(y)) = 0 \text{ for every } y \in T^i \setminus Z\}.$$

In words, $d \in D^i$ assigns numbers to every node in i 's decision tree, with the constraint that zero is assigned to the nodes selected by g^i (one in every set of immediate successors to some nonterminal node.) Note that $\dim(D^i) = \#Z$. Recall that $U^i = R^Z$ denotes the space of payoffs for player i .

For each $(\mu, \pi) \in \Psi$ and $u^i \in U^i$, recursively define a function $v^i(\mu, \pi, u^i)$ on T^i by setting $v^i(\mu, \pi, u^i)(z) = u^i(z)$ for $z \in Z$, and, for $y \in T^i \setminus Z$,

$$(A.3) \quad v^i(\mu, \pi, u^i)(y) = \sum_{y' \in S^i(y)} v^i(\mu, \pi)(y') \cdot v^i(\mu, \pi, u^i)(y').$$

Also define a map $\gamma^i: \Psi \times U^i \rightarrow D^i$ by

$$(A.4) \quad (a) \quad \gamma^i(\mu, \pi, u^i)(y) = v^i(\mu, \pi, u^i)(y) - v^i(\mu, \pi, u^i)(g^i(p^i(y))) \quad \text{for } y \in T^i \setminus \{0\},$$

and

$$(b) \quad \gamma^i(\mu, \pi, u^i)(0) = v^i(\mu, \pi, u^i)(0).$$

In words, we roll-back the tree, obtaining the “value functions” v^i , and then we find the differences in values among branches leading out of each node.

LEMMA A5: *Fixing $(\mu, \pi) \in \Psi$, γ^i is a bijection from U^i to D^i .*

PROOF: We only sketch the proof, as it is straightforward to verify but tedious to write out in complete detail. For fixed (μ, π) and $d^i \in D^i$, we will construct the unique $u^i \in U^i$ such that $\gamma^i(\mu, \pi, u^i) = d^i$. We do this by “rolling forward” through the tree, constructing the corresponding $v^i(\mu, \pi, u^i)$ recursively.

Begin by setting $v^i(\mu, \pi, u^i)(0) = d^i(0)$. Note that if (A.3) and (A.4)(b) are both to be satisfied, we must have

$$\begin{aligned} v^i(\mu, \pi, u^i)(0) &= \sum_{y \in S^i(0)} v^i(\mu, \pi)(y) \cdot [v^i(\mu, \pi, u^i)(g^i(0)) + d^i(y)] \\ &= v^i(\mu, \pi, u^i)(g^i(0)) + \sum_{y \in S^i(0)} v^i(\mu, \pi)(y) \cdot d^i(y), \end{aligned}$$

or, for each $y' \in S^i(0)$,

$$v^i(\mu, \pi, u^i)(y') = v^i(\mu, \pi, u^i)(0) - \sum_{y \in S^i(0)} v^i(\mu, \pi)(y) \cdot d^i(y) + d^i(y').$$

We can repeat this procedure throughout the tree—if we know $v^i(\mu, \pi, u^i)(y)$ for any node y , then (A.3) and (A.4)(b) uniquely determine $v^i(\mu, \pi, u^i)(y')$ for all $y' \in S^i(y)$. Note well that because the transition probabilities out of any node sum to one, there is always a solution to these equations that involves only multiplication, addition, and subtraction. When $v^i(\mu, \pi, u^i)$ has been computed for every y , we have, of course, that the unique corresponding u^i is $u^i(z) = v^i(\mu, \pi, u^i)(z)$. *Q.E.D.*

Define $\phi^i: \Psi \times D^i \rightarrow U^i$ to be the “inverse” of γ^i given in the proof above. Let $D = \prod_{i \in I} D^i$ and $U = \prod_{i \in I} U^i$. Let ϕ be the vector map $(\phi^i)_{i \in I}$ from $D \times \Psi$ to U . Since ϕ consists of iterated multiplications, additions and subtractions (no divisions!), we have the following result.

LEMMA A6: *The map ϕ is smooth (infinitely differentiable in all of its arguments). In fact, ϕ can be extended smoothly to an open domain containing its domain of definition.*

The reason for this construction can now be stated. Fix any consistent basis b , and let $\beta \subseteq A$ be such that $\beta \supseteq A \cap b$. Let D_β be the subset of D where, for each h , if $u(h) = i$, then $d'(a) = d'(a')$ for all $a, a' \in A(h) \cap \beta$ and $d'(a) > d'(a')$ for all $a \in A(h) \cap \beta$ and $a' \in A(h) \cap (A \setminus \beta)$. Note that D_β is a manifold of dimension $\# \beta - \# H$ less than the dimension of D .

LEMMA A7: *The assessment $(\mu, \pi) \in \Psi$ is in $\Phi_b(u)$ if and only if, for some $\beta \supseteq b \cap A$, $(\mu, \pi) \in \text{proj}_{\Psi_b}(\phi^{-1}(u) \cap [D_\beta \times \Psi_b])$. Moreover, $\Phi_b^S(u) = \text{proj}_{\Psi_b}(\phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b])$.*

In words, if we look at ϕ on the domain $\bigcup_{\beta: \beta \supseteq b \cap A} D_\beta \times \Psi_b$, then the set of pre-images of any $u \in U$ is precisely the set of equilibria for u that belong to the basis b . This lemma is a direct application of the definitions of D and the map ϕ and Lemma A4.

A.3. PROOFS OF THE THEOREMS

PROOF OF THEOREM 1: A point $u \in U$ will be called *nice* if for every consistent basis b and every subset β of A such that $\beta \supseteq b \cap A$, the map ϕ restricted to the domain $D_\beta \times \Psi_b$ has u as a regular value. If u is a critical value for any of these domains, u will be said to be *not nice*.

Because the number of pairs (β, b) as above is finite, Sard's Theorem (see Milnor [9]) implies that the set of not nice u has Lebesgue measure zero. We also assert that the set of not nice u is closed. Suppose $\{u_n\}$ is a sequence of not nice points that converges to some u . Let (d_n, μ_n, π_n) be a pre-image of u_n at which the map is critical. We can bound the d_n by the bounds on the μ_n, π_n , and u_n , so without loss of generality we can assume that the (d_n, μ_n, π_n) lies in some single $D_\beta \times \Psi_b$ and that they converge to some (d, μ, π) . If this limit (d, μ, π) lies within the manifold $D_\beta \times \Psi_b$, then continuity of the Jacobian of ϕ on this domain assures us that u must also be a critical value, and hence not nice. If (d, μ, π) lies on the frontier of $D_\beta \times \Psi_b$, then greater care must be taken. In this case (d, μ, π) lies in some $D_{\beta'} \times \Psi_{b'}$ where $\beta' \supseteq \beta$ and $b' \subseteq b$. One can show (using the first part of the Appendix) that any vector in the tangent map at (d, μ, π) in $D_{\beta'} \times \Psi_{b'}$ can be approached by vectors in the tangent maps at the (d_n, μ_n, π_n) . Since the Jacobian is continuous, the criticality of the points (d_n, μ_n, π_n) implies that (d, μ, π) is a critical point, and thus that u is not nice.

Summing up so far, the set of not nice u is a closed set of measure zero. The set of nice u is generic.

Fix any nice u and basis b . From Lemma A7,

$$\Phi_b^S(u) = \text{proj}_{\Psi_b}(\phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b]).$$

Because u is nice, a standard application of the regularity of u implies that $\phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b]$ is either empty or is a manifold of dimension

$$\begin{aligned} \dim(D_{b \cap A}) + \dim(\Psi_b) - \dim(U) &= \#Z \cdot \#I - (\#(b \cap A) - \#H) \\ &\quad + \dim(\Psi_b) - \#Z \cdot \#I \\ &= n(b). \end{aligned}$$

The first part of Theorem 1 is completed by noting that if η is a diffeomorphism carrying some neighborhood of $\phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b]$ into $R^{n(b)}$, then $\eta^*(\mu, \pi) = \eta(\gamma(\mu, \pi, u), \mu, \pi)$ is a diffeomorphism carrying the Ψ_b -projection of this neighborhood into $R^{n(b)}$.

For the second part of Theorem 1, note that the sets D_β for $\beta \supseteq A \cap b$, $\beta \neq A \cap b$, form regular boundaries of $D_{b \cap A}$ —they are hyperplane restrictions. Thus an easy extension of Milnor [9, Lemma 4] yields that for nice u ,

$$\phi^{-1}(u) \cap \left(\bigcup_{\substack{\beta: \beta \supseteq b \cap A \\ \beta \neq b \cap A}} [D_\beta \times \phi_b] \right)$$

is precisely the frontier of $\phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b]$ in $D \times \Psi_b$. The same is clearly true when we project onto Ψ_b . Applying Lemma A7 completes the proof. Q.E.D.

PROOF OF THEOREM 2: If $(\mu, \pi) \in \Phi_b(u)$, then $P^\pi(\cdot)$ depends on π restricted to actions a that are basic and that, in some player's decision tree, have only basic predecessors. Call this set of actions $A(b)$. We therefore will have proved the theorem when we show that for generic u , there is for each b a finite number of $\pi | A(b)$ such that $(\mu, \pi) \in \Phi_b(u)$.

Fixing a basis b , consider the extensive form constructed by taking the tree T and deleting all nodes x whose actions $\alpha(x)$ are not in b , together with all successors of such nodes. Call the resulting tree (and, loosely, the resulting extensive form) T_b . Note that in T_b , the set of actions is $A(b)$. If $(\mu, \pi) \in \Psi_b$ is an equilibrium for u in the original extensive form, then (μ, π) restricted to T_b is surely an equilibrium for u in T_b . Moreover, π restricted to T_b is strictly positive.

Call u very nice if it is nice for the original tree and if, moreover, its projection onto relevant endpoints $T_b \cap Z$ is nice for the games given by T_b , for each b . Call u not very nice if this fails. From the proof of Theorem 1, we know that the set of not very nice points is closed and has Lebesgue measure zero. That is, generic u are very nice.

Fix a very nice u . Applying Theorem 1, the set of strict equilibria for the game T_b where π is strictly positive is a manifold of dimension zero. Moreover, there are no weak equilibria for the game with tree T_b wherein all actions in each information set are equally good. (This simply requires applying Theorem 1 and counting dimensions.) So the set of strict equilibria for the game T_b where π is strictly positive must be finite: If the set were infinite, then it would have a limit point, because the bound on μ and π combined with fixed u yields a bound on d . Any limit point is an equilibrium in which all actions in each information set are equally good. And no such limit point can be strictly positive—it then would not be part of a manifold of dimension zero—nor can it have some action taken with zero probability—for very nice u there are no equilibria for the game with tree T_b wherein all actions in each information set are equally good and where some action has probability zero.

Putting everything together, this shows that for generic (very nice) u , Δ_u must be finite. Q.E.D.

Some comments are in order. As an immediate corollary to this, for generic normal form games there are a finite number of equilibria, none of which is weak. Secondly, Theorem 2 as stated applies equally well to the standard Nash equilibrium definition, by exactly the same argument. (If π is a Nash equilibrium for T , and $A(b)$ is formed in the fashion above, then π must be a strictly positive Nash equilibrium (hence sequential) for T_b .) Thirdly, we can sharpen this result and show that in neighborhoods of generic u , the correspondence $u \Rightarrow \Delta_u$ consists of a finite number of differentiable functions.

Concerning this last comment, we note that it ought to be possible to sharpen Theorem 1 to something like the following. In neighborhoods of generic u , the correspondences $u \Rightarrow \Phi_b^S(u)$ (for $b \in B$) are smooth deformations of manifolds.

PROOF OF THEOREM 3: Call u wonderful if u is very nice and if, moreover, u is a regular value of the mapping ϕ on each of the manifolds $D_\beta \times \Psi_b^0$ for each $b \in B$ and $\beta \supseteq b \cap A$ (cf. Lemma A3). The argument in the proof of Theorem 1 easily extends to show that the set of wonderful u is generic.

Fix a wonderful u , and let (μ, π) be a strict sequential equilibrium. In particular, let $(\mu, \pi) \in \Phi_b^S(u)$. Then there exists $d \in D_{b \cap A}$ such that $(d, \mu, \pi) \in \phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b]$. By Milnor [9, Lemma 4], there exists a sequence $\{(d_n, \mu_n, \pi_n)\} \subseteq \phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b^0]$ with limit (d, μ, π) . Since each $(d_n, \mu_n, \pi_n) \in \phi^{-1}(u) \cap [D_{b \cap A} \times \Psi_b^0]$, each (μ_n, π_n) is a "constrained equilibrium," where the constraints for (μ_n, π_n) are $\pi(a) \geq \min(1/n, \pi_n(a))$ if $a \in b$ and $\pi(a) \geq \pi_n(a)$ if $a \notin b$. These constraints vanish as $n \rightarrow \infty$, so by Selten's first definition of a perfect equilibrium, (μ, π) is perfect.

To obtain the second half of Theorem 3, again fix a wonderful u and any sequential equilibrium (μ, π) . By Theorem 1, we know that there are strict sequential equilibria arbitrarily close to (μ, π) . Moreover, we claim that if (μ', π') is a strict sequential equilibrium close enough to (μ, π) , we have $P^\pi(\cdot) \equiv P^{\pi'}(\cdot)$. This follows directly from the fact that Δ_u is finite. But by the above, (μ', π') is perfect, so that $P^\pi(\cdot)$ is in the projection of the set of perfect equilibria onto Δ . Q.E.D.

REFERENCES

[1] DEBREU, G.: "Economies with a Finite Set of Equilibria," *Econometrica*, 38(1970), 387-392.
 [2] FUDENBERG, D., AND J. TIROLE: "Sequential Bargaining with Incomplete Information," mimeo, Massachusetts Institute of Technology, 1981.
 [3] HARSANYI, J.: "Games with Incomplete Information Played by Bayesian Players," Parts I, II, and III, *Management Science*, 14(1967-68), 159-182, 320-334, 486-502.

- [4] ———: "Advances in Understanding Rational Behavior," Paper CP-366, Center for Research in Management Science, University of California, Berkeley, 1975.
- [5] KOHLBERG, E.: Private Communication, 1980.
- [6] KREPS, D., AND R. WILSON: "Reputation and Imperfect Information," draft, Stanford University, 1981, forthcoming in the *Journal of Economic Theory*.
- [7] KUHN, H.: "Extensive Games and the Problem of Information," in *Contributions to the Theory of Games*, Vol. 2, ed. by H. Kuhn and A. Tucker. Princeton: Princeton University Press, 1953, pp. 193–216.
- [8] MILGROM, P., AND J. ROBERTS: "Limit Pricing and Entry Under Incomplete Information: An Equilibrium Analysis," *Econometrica*, 50(1982), 443–460.
- [9] MILNOR, J.: *Topology from the Differentiable Viewpoint*. Charlottesville, Va.: University Press of Virginia, 1965.
- [10] MYERSON, R.: "Refinements of the Nash Equilibrium Concept," *International Journal of Game Theory*, 7(1978), 73–80.
- [11] NASH, J.: "Noncooperative Games," *Annals of Mathematics*, 54(1951), 286–295.
- [12] RUBENSTEIN, A.: "Perfect Equilibrium in a Bargaining Model," mimeo, Nuffield College, Oxford, 1980.
- [13] SAVAGE, L. J.: *The Foundations of Statistics*. New York: John Wiley and Sons, 1954.
- [14] SELTEN R.: "Spieltheoretische Behandlung eines Oligopolmodells mit Nachfrageträgheit," *Zeitschrift für die gesamte Staatswissenschaft*, 121(1965), 301–324.
- [15] ———: "Re-examination of the Perfectness Concept for Equilibrium Points in Extensive Games," *International Journal of Game Theory*, 4(1975), 25–55.
- [16] WILSON, R.: "Computing Equilibria of Two-person Games from the Extensive Form," *Management Science*, 18(1972), 448–460.