# 20 Acceptable Points in General Cooperative *n*-Person Games

## 1 Introduction

In this paper, we introduce a new theory for *n*-person games. Like many previous theories for *n*-person games, it is based on the search for a reasonable "steady state" in the "supergame,"[1] i.e., the game each play of which consists of a long sequence of plays of the original game. We use two approaches: The usual approach, in which a single play is studied, and an attempt is made to find an *n*-tuple of strategies that intuitively speaking might be acceptable for a steady state in the supergame (Sections 4 and 5); and an approach based on a mathematical analysis of the supergame itself (Sections 6 and 7). The results we obtain by the two approaches turn out to be essentially equivalent (Sections 9 and 10).

We will consider here only cooperative games, i.e., games in which coalitions are permitted; side payments will be forbidden.[2] However, the ideas introduced here also apply to non-cooperative games. Subsequent papers will consider non-cooperative as well as cooperative games.

We will not explicitly consider games involving chance; however, everything we say here carries over with no essential change to games involving chance.

Section 2 describes the notation we will use throughout. The reader may omit this section at first and use it only for reference. Section 8 is devoted to the proof of lemmas needed in the subsequent sections; it makes strong use of the approachability-excludability theory of Blackwell [4]. Sections 11, 12, and 13 are devoted to some miscellaneous remarks and counterexamples.

## 2 Notation

Throughout the paper, we will be concerned with an *n*-person game $G$. $N$ will denote the set of all players; the individual players will be denoted by the positive integers $1, \ldots, n$. If $x^i$ is an object defined for each $i \in N$, and if $B \subset N$, then $x^B$ will denote the set of $x^i$ where $i \in B$; it will be called a $B$-vector. When $B = N$, we will often omit the superscript and the prefix; that is, we will write $x$ instead of $x^N$, and will refer to $x$ simply as a vector rather than as an $N$-vector. If $X$ is a set of vectors, we denote by $X^B$ the

1. The name is due to Luce and Raiffa [6].

2. Of course it is known (see [1]) that the cooperative game with side payments is a special case of the cooperative game without side payments.

set of all $x^B$ where $x \in X$. When we speak of a subset $B$ of $N$, i.e., of a $B \subset N$, we will mean a non-null subset, unless we specifically include the null subset. All $B$-vector equalities or inequalities will be taken to hold term by term.

We will be concerned only with random variables taking a fixed finite number of values, and possibly with infinite sequences of such random variables. If $x$ is a random variable ranging over a finite set $Z$, then the probability distribution of $x$ may be considered as a function $y$ from $Z$ to the reals satisfying the relations:

$$y(z) \geqslant 0, \quad z \in Z \tag{2.1}$$

and

$$\sum_{z \in Z} y(z) = 1. \tag{2.2}$$

We can then write

$$y = \sum_{z \in Z}^{*} y(z)z \tag{2.3}$$

where the $*$ indicates that the sum is merely symbolic and is not meant to be an ordinary sum. $y(z)$ is simply the probability that a random variable distributed according to $y$ will take the value $z$.

The set of all probability distributions on $Z$ will be denoted by $C(Z)$. The set of all pure strategies $p^i$ for player $i$ will be denoted by $P^i$. $P$ is then the space of pure strategy vectors $p$. $C(P)$ is the space of *correlated* strategy vectors. A correlated strategy vector is a probability distribution on the set of all pure strategy *vectors*, just as a mixed strategy is a probability distribution on individual strategies. In general, to make use of a correlated strategy vector, the players have to agree to consult the *same* random device in choosing the pure strategies they will play. Correlated strategy vectors are thus not usable in non-cooperative games. We will usually abbreviate $C(P)$ to $C$. The prefix $c$- in front of a word will stand for "correlated."

Let $X^B$ be a space of $B$-vectors, and let $B = \bigcup_{j=1}^{k} B_j$ be a partition of $B$ into subsets. If for $j = 1, \ldots, k$, we have

$$y^{B_j} \in C(X^{B_j}),$$

then $(y^{B_1}, \ldots, y^{B_k})$ denotes the joint probability distribution of $y^{B_1}, \ldots, y^{B_k}$.

If $Z_1$ and $Z_2$ are finite sets and $y \in C(Z_1 \times Z_2)$, then $y|Z_1$ (sometimes also written $y|C(Z_1)$) will be the distribution given by

$$y|Z_1 = \sum_{z_1 \in Z_1}^{*} z_1 \sum_{z_2 \in Z_2} y(z_1, z_2).$$

If $B_1 \subset B_2$ and $y^{B_2} \in C(X^{B_2})$, then we will simply write $y^{B_1}$ instead of $y^{B_2}|X^{B_1}$. This applies to all symbols *except* $\gamma$. Thus $\gamma^{B_1}$ is *not* the same as $\gamma^{B_2}|P^{B_1}$. ($\gamma^B$ will be defined on $C^B(P)$.)

Let $Z_1$ and $Z_2$ be finite sets, and suppose we have a function $f$ that takes $Z_1$ into $Z_2$. If nothing is said to the contrary, we will assume that the definition of $f$ is extended linearly; that is, $f$ is a function from $C(Z_1)$ to $C(Z_2)$ defined by

$$f(y) = \sum_{z_1 \in Z_1}^{*} f(z_1)y(z_1).$$

If for each $i \in N$, we have a function $f^i$ defined on $X^i$, then we may define a *vector function f* on $X$ by means of

$$f(x) = (f^1(x^1), \ldots, f^n(x^n)) \tag{2.4}$$

for all $x \in X$. It is important to remember, though, that the inverse process does not in general work. A vector function $f$ defined on $X$ does not yield a unique definition of a function $f^i$ defined on $X^i$.

In general, when we have defined $f^i$ on $X^i$, we will use the function $f$ on $X$ in the above sense.

For $p \in P$, the payoff to the $i^{\text{th}}$ player will be denoted by $E^i(p)$. If $c \in C$, $H(c)$ will be the mean of the distribution $E(c)$. $E(c)$ is called the *payoff distribution* when $c$ is played, $H(c)$ the *expected* payoff when $c$ is played. If $p \in P$, we have $H(p) = E(p)$.

The null set will be denoted by $\varnothing$.

We will be concerned with partitions of $N$. Let $\mathscr{P}_N$ denote the set of all subsets of $N$. Formally, a partition of $N$ is a vector $d$ whose components are members of $\mathscr{P}_N$, and which satisfies the condition

$$j \in d^i \Longleftrightarrow d^i = d^j, \quad i, j \in N.$$

The set of all partitions of $N$ will be denoted by $D$. The *constant* partition $d_N$ is defined by

$$d_N^i = N, \quad i \in N. \tag{2.5}$$

The *identity* partition $d_e$ is defined by

$$d_e^i = \{i\}, \quad i \in N.$$

Denote by $R$ the set of all vectors $r$, each of whose components $r^i$ is a member of $\mathscr{P}_N$ containing $i$. Clearly $D \subset R$, and $D^i = R^i$. For a given $d$, the number of distinct $d^i$ will be called $n(d)$. We choose a set of representatives $i_1, \ldots, i_{n(d)}$, one out of each distinct $d^i$. If $c^B$ is a $c$-strategy $B$-vector, define

$e(c^B) = B.$

(The $e$ stands for exponent; $B$ is called the *exponent* of $c^B$.) Denote by $T$ the set of all vectors $t$, each of whose components is some correlated strategy $B$-vector (not necessarily the same $B$ for each member of $N$), for which

$i \in e(t^i), \quad i \in N.$

Let

$$d^i(t) = \{j | t^j = t^i\}, \quad i \in N, \quad t \in T. \tag{2.6}$$

Then $d(t)$ is a partition of $N$, and we have

$d^i(t) \subset e(t^i), \quad i \in N.$

Define

$$c(t) = ((t^{i_1})^{d^{i_1}(t)}, \ldots, (t^{i_{n(d(t))}})^{d^{i_{n(d(t))}}(t)}) \tag{2.7}$$

where $i_1, \ldots, i_{n(d(t))}$ are a set of representatives for $d(t)$.

In discussing games in extensive form, we will call the vertices of the game tree that are not moves by the name "terminal" rather than "play." This is necessary to avoid confusion, as we will be dealing with repeated plays of the same game, in the more common sense of the word.

As a final remark, we mention that everything we have stated probabilistically can be restated, where necessary, in the language of measure theory, and rigorously justified in that language. All sets we use can be proved to be measurable.

## 3    General Cooperative Games

We use the word "general" in the title of the paper and of this section in order to emphasize that, unlike the cooperative games of von Neumann–Morgenstern, the games to be considered here exclude side payments. On the other hand, full and free cooperation is permitted between the players, in the sense that they may communicate with each other freely prior to each play and form any coalitions they may consider convenient. The agreements under which the coalitions are chosen are to be considered enforceable. (However, see Section 13 for further discussion of this point.) We will usually omit the word "general" and refer simply to a cooperative game. This is in accordance with the usage introduced by Nash [1].

In applications, the prohibition on side payments may be the result either of a physical barrier to side payments (in economic applications this might take the form of certain provisions in an anti-trust law) or of the lack of a common unit of measurement for the payoff. As has been pointed out by Raiffa [8], if there is a common unit of measurement for the payoff but there exists some physical barrier to the making of side payments, then the theory need only be invariant under linear transformations of the payoff function that do not destroy the common unit of measurement; i.e., under transformations of the form

$$H^! = aH + b \tag{3.1}$$

where both $a$ and $b$ are scalar multiples of the vector $(1, \ldots, 1)$ ($a$ by a positive factor, $b$ by an arbitrary factor). If, however, there is no common unit of measurement for the payoff, then the theory should be invariant under independent linear transformations of the payoff function; i.e., under transformations of the form (3.1), where $a$ is an arbitrary positive vector and $b$ is an arbitrary vector. In fact, the theory presented in this paper is invariant under the wider class of transformations, so that the results may be used in either kind of application.

Mathematically, the cooperative game is given by a vector payoff function that is assumed to be linear on the space $C$ of correlated strategy vectors. The cooperative game differs from the non-cooperative game only in that the use of correlated strategy vectors is permitted in the former, but not in the latter.

This points up a reason for studying cooperative games even if our a priori interest lies in the field of non-cooperative games only. Consider a long sequence of plays (we call this a *superplay*) of the game

|   | 1 | 2 |
|---|-------|-------|
| 1 | (4, 0) | (0, 0) |
| 2 | (0, 0) | (0, 4) |

.

One of the most "natural" ways for the players to proceed in such a superplay would be to alternate between playing (1, 1) and (2, 2). This would tend to happen even in the non-cooperative case, by a sort of "silent gentlemen's agreement." However, the resulting payoff, (2, 2), is obtainable only by a correlated strategy vector, not by mixed strategies. Thus if we wish to arrive at results that are valid for the supergame by considering a single play, we must permit ourselves the use of correlated strategy vectors. This amounts to considering cooperative games.

# 4  Acceptable Points

Let $c$ be a correlated strategy vector in $G$. Obviously, $c$ will be a candidate for a possible steady state in the supergame unless some subset $B$ of the players can, by concerted action, increase their payoff. Here it is not sufficient that the players of $B$ be able to increase their payoffs by changing their strategies while the players of $N-B$ maintain their strategies as they are. For the players of $N-B$ will not in general maintain their strategies fixed in the face of a change on the part of $B$, but will also change *their* strategies in order to meet the new conditions; thus $B$'s glory would be short-lived indeed. In order to rule out $c$ as a candidate for a possible steady state, we have to be sure that the players of $N-B$ cannot prevent the players of $B$ from improving their payoff by concerted action. Of course, we must require that *all* the members of $B$ obtain more at the new point than they do at $c$, for they must all have an incentive to cooperate in the venture.

DEFINITION 4.1    Let $c_0 \in C$. $c_0$ is said to be $c$-acceptable if there is no $B \subset N$ such that for each $c^{N-B} \in C^{N-B}$, there is a $c^B \in C^B$ for which

$$H^B(c^B, c^{N-B}) > H^B(c_0).$$

We write "$c$-acceptable" instead of "acceptable" because in subsequent papers of this series, we will be defining a slightly different kind of acceptability, and we will use prefixes to distinguish them. When in heuristic discussion in the sequel we use the word "acceptable" without a prefix, then it is to be taken to mean "acceptable" with an arbitrary prefix (but one that is fixed throughout the discussion). For the purposes of this paper alone, "acceptable" can be taken to mean the same thing as "$c$-acceptable."

The set of all $c$-acceptable $c$-strategy vectors is denoted by $A_c$. Note that the notion of acceptability is a "global" one. In other words, whether or not a point is acceptable depends only on the payoff at that point, not on the point itself. This is also easy to accept from the intuitive standpoint; as long as he is getting the same payoff, it makes little difference to a player which strategy he is playing.

DEFINITION 4.2    A payoff vector $h$ is said to be $c$-acceptable, if for some $c \in A_c$, we have

$$H(c) = h.$$

It will be convenient in the sequel to have available the following trivial restatement of 4.2.

LEMMA 4.3   A payoff vector $h$ is $c$-acceptable if and only if for each $B \subset N$, there is a $c^{N-B} \in C^{N-B}$, such that for all $c^B \in C^B$, there is an $i \in B$ for which

$$H^i(c^B, c^{N-B}) \geqslant h^i.$$

## 5   Two-Person, Zero-Sum Games

A sine qua non of theories of $n$-person games is that for two-person, zero-sum games, they reduce to the von Neumann theory. Because of the global nature of our theory, we can only expect that every acceptable payoff vector yield to each player the value of the game to that player. This is indeed the case.

THEOREM 1   In a two-person, zero-sum game $G$, a payoff vector $h$ is $c$-acceptable if and only if

$$h^1 = v, \tag{5.1}$$

where $v$ is the value of $G$ to player 1.

*Proof*   Let $h$ be $c$-acceptable. From Lemma 4.3 applied to $B = \{1\}$, we deduce that there is a $c_0^2 \in C^2$, such that for all $c^1 \in C^1$, we have

$$H^1(c^1, c_0^2) \leqslant h^1.$$

Hence

$$\max_{c^1 \in C^1} \; H^1(c^1, c_0^2) \leqslant h^1.$$

Hence

$$\min_{c^2 \in C^2} \; \max_{c^1 \in C^1} \; H^1(c^1, c^2) \leqslant h^1. \tag{5.2}$$

$C^1$ and $C^2$ are the mixed strategy spaces for players 1 and 2 respectively. Hence the left side of (5.2) is $v$, and we deduce

$$v \leqslant h^1. \tag{5.3}$$

The value of $G$ to player 2 is $-v$. Proceeding as in the proof of (5.3), we obtain

$$-v \leqslant h^2. \tag{5.4}$$

But

$$h^2 = -h^1. \tag{5.5}$$

Combining (5.4) and (5.5), we obtain

$$v \geqslant h^1. \tag{5.6}$$

Combining (5.3) and (5.6), we obtain (5.1).

Conversely, assume (5.1). Then (5.3) and (5.2) follow at once. But for $B = \{1\}$, (5.2) is precisely the condition for $c$-acceptability given by Lemma 4.3. Similarly we establish the condition for $B = \{2\}$. To establish it for $B = \{1, 2\}$, note that for *any* $(c^1, c^2)$, we have

$$H^1(c^1, c^2) + H^2(c^1, c^2) = 0 = h^1 + h^2,$$

whence we must have either

$$H^1(c^1, c^2) \leqslant h^1$$

or

$$H^2(c^1, c^2) \leqslant h^2.$$

Thus $h$ is $c$-acceptable.

## 6  Strategies and Payoffs in the Supergame

In the notion of supergame that will be used in this paper, each superplay consists of an *infinite* number of plays of the original game $G$. On the face of it, this would seem to be unrealistic, but actually it is more realistic than the notion in which each superplay consists of a fixed finite (large) number of plays of $G$. In the latter notion, the fact that the players *know* when they have arrived at the last play becomes the decisive factor in the analysis, overshadowing all other considerations.[3] This is unnatural, because in the usual case, the players can always anticipate that there will still take place in the future an indefinite number of plays. This is especially true in political and economic applications, and even holds true for the neighborhood poker "club." Of course when looked at in the large, nobody really expects an infinite number of plays to take place; on the other hand, after each play we do expect that there will be more. A. W. Tucker has pointed out that this condition is mathematically equivalent to an infinite sequence of plays, so that is what our notion of supergame will consist of.

Roughly speaking, we will be interested in "behavior" strategies in the supergame; that is, we will be interested in strategy vectors that determine

---

3. See, for example, the excellent analysis of the supergame of the game known as the "Prisoner's Dilemma" that appears in 5.5 of [6].

the *c*-strategy vector to be played on each play as a function of the information possessed by the players about the outcomes of the previous plays. Before we pass to the formal definition of a supergame *c*-strategy vector, let us examine in detail the procedure followed during the play of a cooperative supergame. The first question to be answered is: What information about the outcome of each previous play does each player $i$ actually possess? The most he can know is exactly what pure strategy vector was actually played; the least he can be sure of knowing, in every game, is the strategy he himself played and the payoff he himself received. In between is a whole range of possibilities, depending on the rules of *G*. The situation can be concisely described by means of a function $u^i$ that takes $P$ onto some identification space of $P$. We define $u^i$ so that for $p_1, p_2 \in P$, $u^i(p_1) = u^i(p_2)$ if and only if player $i$ cannot tell after a play of *G*, whether the vector $p_1$ or the vector $p_2$ has been played. The minimum knowledge that each player $i$ has, in accordance with our remarks above, insures us that if

$$u^i(p_1) = u^i(p_2),$$

then

$$p_1^i = p_2^i$$

and

$$H^i(p_1) = H^i(p_2).$$

Every game *G* has associated with it an information vector $u$ satisfying these conditions. In particular, these conditions assure us of the existence of functions $\phi^i$ and $\psi^i$ which map $u^i(P)$ onto $P^i$ and $H^i(P)$ respectively. Some relations connecting the vector functions $H$, $\phi$, $\psi$, and $u$ are the following:

$$\phi \circ u = \text{identity}$$

$$\psi \circ u = H$$

$$\psi = H \circ \phi.$$

The information vector just defined is not included in the classical formulation of the extensive form of a game, as given, say, in [3]. That is because the state of information at the *end* of a game has heretofore been of little interest, as the payoff in no way depends on it. In a number of consecutive plays of a game, however, it is seen to be of vital importance. Incidentally, inclusion of this concept in the extensive structure of a game necessitates a revision in the notion of equivalence as formulated in [5].

We hope to treat the subject in more detail in a subsequent paper. For the present, we content ourselves with the remark that the information vector may be included in the extensive description of the game by associating with each player $i$, an information partition on the set of all terminals. In a game of perfect information, this partition would be the identity partition for each $i$.

All that we have said in this section up to now is quite general, applying to the non-cooperative as well as to the cooperative case. In the latter case, the information function does not exhaust the information that a player has to work with. Prior to each play, consultation takes place among the various players to determine the constitution of the coalitions and the correlated strategy $B$-vector that each coalition $B$ will use. Formally, each player designates the coalition $B$ to which he wishes to belong, and the correlated strategy $B$-vector he wishes $B$ to play. Those players who can agree among each other form coalitions and play the agreed upon $c$-strategy $B$-vector.[4]

Naturally, although each player knows nothing about the coalitions in which the other players participate, he does remember the coalition $B$ of which he himself was a member on previous plays.

The formal definition is as follows:

DEFINITION 6.1   A supergame $c$-strategy $f^i$ for player $i$ is an infinite sequence of functions

$$f_0^i, f_1^i, \ldots, f_k^i, \ldots$$

where $f_k^i$ is a function from

$$((U_1^i \times R_1^i) \times \cdots \times (U_k^i \times R_k^i))$$

to $T^i$; here $U_j^i$ is a copy of $u^i(P)$ and $R_j^i$ is a copy of $R^i$, for $j = 1, \ldots, k$.

The vector $f$ is called a supergame $c$-strategy vector.

In the sequel, we will denote $U_j^i \times R_j^i$ by $J_j^i$. $J_j^i$ represents the information about the $j^{\text{th}}$ play available to player $i$.

To define the payoff to a supergame $c$-strategy vector, we must first define the payoff to a given super*play*, as a function of its individual constituent plays. For this purpose, we use the arithmetic mean of the payoffs (in the limiting case, this amounts to the first Cesàro sum); Blackwell [4] also used the arithmetic mean under similar circumstances. Heuristically, one can argue both for and against this method; in my mind the heuristic advantages outweigh the disadvantages. Other methods have been proposed (see [6]) for obtaining the superplay payoff from the pay-

---

4. If $B_1 \subset B_2$ and every member of $B_1$ proposed $(B_2, c^{B_2})$ while the members of $B_2 - B_1$ proposed something else, then the members of $B_1$ form a coalition and play $c^{B_1}$.

offs to its constituent plays; the most prominent seems to be the discounting of future payoffs back to the present, using a fixed rate of interest. This has the disadvantage of putting an unnaturally heavy stress on the beginning of the supergame. Nevertheless, it would be interesting to see an analysis based on this payoff notion, especially for certain kinds of economic applications.

A more perplexing problem in defining the payoff to a supergame *c*-strategy vector $f$ is that of somehow combining the various possible superplays that might occur if $f$ is played. The first impulse is to use some notion of *expected* payoff. We shall indeed define such a notion, but it is not an appropriate criterion for the players in choosing their supergame *c*-strategy. The chief reason for our being interested in expected payoffs for *individual* plays is that in a long sequence of plays, the law of large numbers turns an expected payoff into an *actual* payoff. A payoff for the whole superplay that is only *expected*, but not eventually approached, will not satisfy most players. Thus if a player has joined a coalition on the strength of a high expectation and is being disappointed in this expectation, he will soon express his disappointment by resigning from the coalition, regardless of his originally high expectation. For a clear example of the pitfalls of expected payoffs, see Section 12. What is required is some kind of criterion based on a law of large numbers. We give a restricted definition of this kind in this section, but reasoning of this type in general will dominate the remainder of the paper.

Let us first return to the question of expected payoff. For each $t \in T$, define

$$s(t) = (u(c(t)), d(t)). \tag{6.2}$$

Let $f$ be a supergame *c*-strategy vector. For each $k \geqslant 1$, we may define a member $z_k(f) = z_k$ of $C(J_1 \times \ldots \times J_k)$ recursively as follows:

$$z_1 = s(f_0) \tag{6.3}$$

$$z_k = \sum_{y \in J_1 \times \cdots \times J_{k-1}}^{*} z_{k-1}(y)(y, s(f_{k-1}(y))), \quad k > 1. \tag{6.4}$$

$z_k$ is the probability distribution of possible outcomes of plays up to the $k^{\text{th}}$. Define a member $x_k(f) = x_k$ of $C(J_k)$ as follows:

$$x_k = s(f_{k-1}(z_{k-1})), \quad k \geqslant 1. \tag{6.5}$$

$x_k$ is the probability distribution of outcomes of the $k^{\text{th}}$ play. Finally, define $E_k = E_k(f)$ and $H_k = H_k(f)$ by

$$E_k = E(c(f_{k-1}(z_{k-1}))) \tag{6.6}$$

and

$$H_k = H(c(f_{k-1}(z_{k-1}))). \tag{6.7}$$

$E_k$ represents the probability distribution of possible payoff vectors on the $k^{\text{th}}$ play; $H_k$ is the expected payoff for the $k^{\text{th}}$ play.

If

$$H(f) = \lim_{k \to \infty} \frac{1}{k} \sum_{r=1}^{k} H_r(f) \tag{6.8}$$

exists, then $f$ is said to be *summable in the mean*.

Let $v = (v_1, \ldots, v_k, \ldots)$ be a sequence of random variables in which $v_k$ is distributed according to the distribution $x_k(f)$. These random variables are not in general independent. Instead of saying that the random variables $v_k$ are distributed according to $x_k(f)$, we will sometimes say that the random variable $v$ is distributed according to $f$. The probability of a statement concerning a random variable $v$ distributed according to $f$ will usually be denoted $\text{Prob}_f$.

DEFINITION 6.9    A supergame $c$-strategy vector $f$ is said to be *summable* if it is summable in the mean and if a sequence of random variables distributed according to $E_k(f)$ obeys the strong law of large numbers.[5]

If $v$ is an infinite sequence $(v_1, \ldots, v_k, \ldots)$ where $v_k \in J_k$ for all $k$, define

$$H_k(v) = \psi(v_k|U_k)(= H(\phi(v_k|U_k))) \tag{6.10}$$

and

$$S_k(v) = \frac{1}{k} \sum_{j=1}^{k} H_j(v). \tag{6.11}$$

Note that if $v$ is a random variable distributed according to $f$, then $H_k(v)$ is a random variable distributed according to $E_k(f)$.

LEMMA 6.12    A necessary and sufficient condition that a supergame $c$-strategy vector $f$ be summable is that there exist a vector $H(f)$ such that for each vector $\varepsilon > 0$, we have

$$\lim_{k \to \infty} \text{Prob}_f(|S_r(v) - H(f)| \geqslant \varepsilon \text{ for some } r \geqslant k) = 0.$$

Furthermore, this $H(f)$ is then identical with the $H(f)$ defined by (6.8).

The proof, which is straightforward, will be omitted.

---

5. See [2], p. 207. The law is there stated for one dimensional random variables; the extension to $n$ dimensions is straightforward.

## 7   Strong Equilibrium Points

When we come to investigate what supergame $c$-strategy vectors $f$ are liable actually to occur, we are faced with a problem somewhat different from the one we discussed in Section 4. There one could always look forward to future plays in order to punish deviations by the other players. In the supergame, however, there is only one superplay. Thus the players will have to make sure the first time that deviations won't pay. This amounts to saying that no coalition $B$ can increase the payoff of all its members by substituting a different supergame $c$-strategy $B$-vector $g^B$ while $N - B$ *maintains* the supergame $c$-strategy $(N - B)$-vector $f^{N-B}$. It is clear that such an $f$ would be very strongly stable, since no *coalition* could have any incentive to move away from it, especially if we demand that $B$ cannot even increase its payoff with positive probability, no matter how small. This is a much stronger version of what corresponds in the supergame to an ordinary equilibrium point.

The formal definition is as follows:

DEFINITION   Let $f$ be a summable supergame $c$-strategy vector. $f$ is a *strong equilibrium c-point* if there is no $B \subset N$ for which there is a supergame $c$-strategy vector $g$ satisfying

$$g^{N-B} = f^{N-B} \tag{7.1}$$

and

$$\lim_{k \to \infty} \mathrm{Prob}_g(S_r^B(v) \geqslant H^B(f) + \varepsilon^B \text{ for some } r \geqslant k) > 0 \tag{7.2}$$

for some $B$-vector $\varepsilon^B > 0$.

(The limit exists because the sequence involved is monotone decreasing.)

The set of all strong equilibrium $c$-points is denoted by $S_c$.

It is necessary to use the form $S_r^B(v) \geqslant H^B(f) + \varepsilon^B$ rather than $S_r^B(v) > H^B(f)$ because substitution of the latter statement would make condition 7.2 true even for $f = g$. We wish to include only those cases in which substitution of $g$ for $f$ will yield $B$ an advantage that at least does not tend to the vanishing point, doubtful and rare as it may be.

Condition 7.2 is very weak, which means that the concept of strong equilibrium $c$-point is a very strong one (i.e., comparatively few supergame $c$-strategy vectors are strong equilibrium $c$-points). On the other end of the spectrum, we could replace 7.2 by

$$\lim_{k \to \infty} \text{Prob}_g(S_r^B(v) \geqslant H^B(f) + \varepsilon^B \text{ for all } r \geqslant k) = 1$$

for some $B$-vector $\varepsilon^B > 0$.    (7.3)

Denote by $\tilde{S}_c$ the set of supergame $c$-strategy vectors obtained when 7.3 is substituted for 7.2 in the definition of strong equilibrium $c$-point.

It is easy to see that $7.3 \Longrightarrow 7.2$, whence it follows that

$$S_c \subset \tilde{S}_c$$

and

$$H(S_c) \subset H(\tilde{S}_c).$$    (7.4)

We will demonstrate the opposite inequality as well. This shows that the notion of strong equilibrium point does not really depend on the strength of the condition 7.2.

## 8    Application of Approachability and Excludability Theory

For two-person, zero-sum games with vector payoffs, Blackwell [4] has defined the notion of "approachability" or "excludability" of a set in the payoff space with a supergame strategy employed by one or the other of the players. Blackwell does not use the information vector $u$ which we introduced in Section 6, but apart from that, our Definition 6.1 applied to two-person, zero-sum games reduces to his definition of supergame strategy (for games not involving chance). Most of his definitions and results can be used here without any changes being necessary due to our use of the information vector.

Let us consider the two-person, zero-sum game with vector payoffs that is obtained from our game $G$ by considering the coalition $B$ to be the first player and the coalition $N - B$ the second player, the payoff to the first player being given by $H^B$.

LEMMA 8.1    Let $\Lambda(f, \varepsilon^B)$ denote the closed convex set consisting of all $B$-vectors $\lambda^B$ for which

$$\lambda^B \geqslant H^B(f) + \varepsilon^B.$$

Then if there is an $\varepsilon^B > 0$ such that $\Lambda(f, \varepsilon^B)$ is approachable with $g^B$, then 7.3 holds.

*Proof*    This lemma follows at once from the definition of approachability; see Section 1 of [4]. Our games do not involve chance, so the elements of the game matrix should be considered as $B$-vectors rather than probability distributions on $B$-vectors. We have abbreviated

"approachable in the game matrix" to "approachable"; this will be done throughout the sequel.

LEMMA 8.2   If there is an $\varepsilon^B > 0$ such that, for each $c^{N-B} \in C^{N-B}$, there is a $c^B \in C^B$ for which

$$H^B(c^B, c^{N-B}) \geqslant H^B(f) + \varepsilon^B,$$

then there is a $g^B$ for which 7.3 holds.

*Proof*   The condition given in the lemma is equivalent to the statement that for each $c^{N-B} \in C^{N-B}$, $H^B(C^B, c^{N-B})$ intersects $\Lambda(f, \varepsilon^B)$. In accordance with Theorem 3 of [4], this in turn is equivalent to the statement that there is a $g^B$ with which $\Lambda(f, \varepsilon^B)$ is approachable. Applying 8.1, we obtain the conclusion.

The information situation in [4] is slightly different from ours. In our application, the requirements of [4] may be stated as follows: Both $B$ and $N - B$ know the payoffs to $B$ on each previous play; nothing else is known. In our situation, $B$ knows its own previous payoffs, but $N - B$ need not know what the payoffs to $B$ were; on the other hand, both $B$ and $N - B$ may have information given by the information function, that they do not have in [4]. An analysis of the proof of Theorem 3 of [4] shows that these differences do not affect the truth of the theorem, so that we may still apply it here.

To see this intuitively, note that the only difference that could interfere with the approachability of $\Lambda(f, \varepsilon^B)$ with $g^B$ would be the additional information that $N - B$ has in our case that it does not have in [4]. Examination of the proof of Theorem 3 of [4] indicates that $g_k^B$ depends only on the previous payoffs, and can be defined for *any* set of previous payoffs without affecting the truth of the theorem; thus no amount of additional information can possibly help $N - B$ in preventing $B$ from approaching $\Lambda(f, \varepsilon^B)$ with $g^B$.

LEMMA 8.3   If there is a $c^{N-B} \in C^{N-B}$ (called $\gamma^{N-B}$) such that for each $c^B \in C^B$, there is an $i \in B$ for which

$$H^i(c^B, \gamma^{N-B}) \leqslant H^i(f),$$

then for every $B$-vector $\varepsilon^B > 0$, and every supergame $c$-strategy vector $\beta$ for which for all $(v_1, \ldots, v_k)$, we have

$$\beta_k^i(v_1, \ldots, v_k) = \gamma^{N-B}, \quad i \in N - B, \quad k \geqslant 0, \tag{8.4}$$

we have

$$\lim_{k \to \infty} \mathrm{Prob}_\beta(S_r^B(v) \geqslant H^B(f) + \varepsilon^B \text{ for some } r \geqslant k) = 0.$$

*Proof*   The hypothesis of the lemma is equivalent to the assertion that for each $\varepsilon^B > 0$, $\Lambda(f, \varepsilon^B)$ fails to intersect $H(C^B, \gamma^{N-B})$. Applying Theorem 3 of [4], we conclude that for each $\varepsilon^B > 0$, $\Lambda(f, \varepsilon^B)$ is excludable with the supergame $c$-strategy $(N - B)$-vector given by (8.4). From the definition of excludability (Section 1 of [4]) it then follows that there is a number $\delta > 0$ such that

$\lim_{k \to \infty} \text{Prob}_\beta$ (For all $r \geqslant k$, there is an $i \in B$ such that

$S_r^i(v) < H^i(f) + \varepsilon^i - \delta) = 1,$

whence

$$\lim_{k \to \infty} \text{Prob}_\beta(S_r^B(v) \geqslant H^B(f) + \varepsilon^B \text{ for some } r \geqslant k) = 0.$$

## 9   The Main Theorem: First Half

We will now show that our two approaches culminating in the definition of $c$-acceptability (4.1) and strong equilibrium (7.1, 7.2, and 7.3) actually yield results that are essentially the same. More precisely, for every $c$-acceptable strategy vector $c$, there is a strong equilibrium $c$-point (according to either 7.2 or 7.3) in supergame $c$-strategies that has the same payoff as $c$; and conversely. Thus $c$-acceptable points for a single play are seen to correspond exactly to strong equilibrium $c$-points in the supergame. The main theorem can be concisely stated by means of the equation

$$H(A_c) = H(S_c)(= H(\tilde{S}_c)).$$

This section is devoted to the proof of the relation

$$H(\tilde{S}_c) \subset H(A_c).$$

The proof rests heavily on Lemma 8.2. In addition to 8.2, the main fact needed is that if $B$ cannot be prevented from obtaining more than $h^B$, then they cannot be prevented from obtaining a fixed amount more than $h^B$ either. This is shown in Lemma 9.1.

LEMMA 9.1   Let $h$ be a vector and $B$ be a subset of $N$. If for each $c^{N-B} \in C^{N-B}$, there is a $c^B \in C^B$ for which

$$H^B(c^B, c^{N-B}) > h^B,$$

then there is a positive $B$-vector $\varepsilon^B$ such that for each $c^{N-B} \in C^{N-B}$, there is a $c^B \in C^B$ for which

$$H^B(c^B, c^{N-B}) \geqslant h^B + \varepsilon^B.$$

*Proof* We first remark that if $G(x, y)$ is a continuous real-valued function of two variables $x$ and $y$ ranging over compact sets $X$ and $Y$ respectively, then the function defined by

$$F(y) = \max_{x \in X} G(x, y)$$

is continuous on $Y$.

Now define a function $F$ on $C^{N-B}$ by

$$F(c^{N-B}) = \max_{c^B \in C^B} \min_{i \in B} (H^i(c^B, c^{N-B}) - h^i).$$

By the above remark, $F$ is continuous. Since $C^{N-B}$ is compact, $F$ takes on its minimum at a point $c_0^{N-B}$ in $C^{N-B}$. By the hypothesis of the lemma,

$$F(c_0^{N-B}) > 0.$$

Setting

$$\varepsilon^i = F(c_0^{N-B})$$

for all $i \in B$, we obtain the conclusion of the lemma.

THEOREM 2    $H(\tilde{S}_c) \subset H(A_c)$.

*Proof* Let $f$ be a summable supergame $c$-strategy vector; set

$$h = H(f).$$

Suppose

$$h \notin H(A_c), \tag{9.2}$$

i.e., suppose there is a $B$ for which the hypothesis of Lemma 9.1 holds. Then the conclusion of Lemma 9.1, which is the same as the hypothesis of Lemma 8.2, holds. Hence we deduce the conclusion of Lemma 8.2, i.e., there is a $g^B$ for which 7.3 holds. Setting

$$g^{N-B} = f^{N-B},$$

we obtain a $B$ and a $g$ obeying 7.1 and 7.3, whence

$$f \notin \tilde{S}_c. \tag{9.3}$$

As we can go through the same argument for any $f$ satisfying $h = H(f)$, we may deduce from (9.3) that

$$h \notin H(\tilde{S}_c). \tag{9.4}$$

We have shown that $(9.2) \Longrightarrow (9.4)$, whence we can deduce that

$$H(\tilde{S}_c) \subset H(A_c),$$

which is what we set out to prove.

## 10   The Main Theorem: Second Half

This section is devoted to the proof of the relation

$$H(A_c) \subset H(S_c);$$

together with Theorem 2, this yields the main theorem. The train of thought of the proof is somewhat as follows:

Suppose $h$ to be a $c$-acceptable payoff. It is easy to set up a supergame $c$-strategy vector $f$ whose payoff is $h$. But we must also incorporate in $f$ a foolproof system for the punishment of any deviators; $f$ must make sure that crime does not pay. The machinery for accomplishing this is at hand; by the definition of $c$-acceptability, for each $B$ there is a $c$-strategy $(N - B)$-vector $\gamma^{N-B}$ whose use prevents $B$ from obtaining more than $h^B$. It remains only to *find* the culprits. To do this, we note that any players who at one time or another deviated from $f$ were at the time of their deviation certainly not included in the same coalition as the "orthodox" players. Since each player knows who was in his coalition on all previous plays, the deviators are thus easily spotted. The set $B_k$ of players who deviated on some play up to the $(k + 1)^{\text{st}}$ is monotone increasing with $k$. Every time it grows, i.e., every time another player joins the ranks of the deviators, the remaining players revise their strategy to punish at least one of the *new* set of deviators. Since $N$ is finite, there must be a set $B$ that includes all the deviators and that is actually attained by $B_k$ after a finite number of plays. The use of the strategy $\gamma^{N-B}$ will then make certain that crime does not pay for at least one of the deviators. The proof is complicated by the fact that $B_k$ and $B$ usually depend on which pure strategies were chosen from among those considered by the correlated strategy vectors used on previous plays.

The strong equilibrium point $f$ just described is one of "unrelenting ferocity"[6] against offenders. It exhibits a zeal for meting out justice that is entirely oblivious to the sometimes dire consequences to oneself or to the other faithful—i.e., those who have not deviated. There are other, more reasonable strong equilibrium points, which give deviating players a chance to return to the fold. These are important in connection with the general non-cooperative game; I will return to them in a subsequent

6. The term is due to Luce and Raiffa [6].

paper on that topic. Here it is simpler and just as efficacious to use the strong equilibrium point of unrelenting ferocity.

THEOREM 3    $H(A_c) \subset H(S_c)$.

*Proof*   Let $h \in H(A_c)$, and suppose $\gamma^N \in A_c$ is such that

$$H(\gamma^N) = h. \tag{10.1}$$

Then by 4.3, for each $B \subset N$ there is a $c^{N-B} \in C^{N-B}$, which we will call $\gamma^{N-B}$, such that for each $c^B \in C^B$, there is an $i \in B$ for which

$$H^i(c^B, \gamma^{N-B}) \leqslant h^i. \tag{10.2}$$

Now define a supergame $c$-strategy vector $f$ as follows:

$$\begin{aligned}
f_o^i &= \gamma^N, \quad i \in N \\
f_k^i(v_1^i, \ldots, v_k^i) &= \gamma^{v_k^i | R^i}, \quad k > 0, \quad i \in N, \quad v_j^i \in J_j^i, \quad j \leqslant k.
\end{aligned} \tag{10.3}$$

We must prove that the payoff to $f$ is $h$ and that $f$ is a strong equilibrium $c$-point. It is convenient to break the proof into a number of lemmas. However, these lemmas will depend on previous assumptions and formulae, and are valid only within the context of this proof. They will be numbered like ordinary formulae used within the proof.

LEMMA 10.4    For $k > 0$ and $i \in N$, we have

$$f_k^i(z_k) = \gamma^N.$$

*Proof*   The proof is by induction on $k$. For $k = 0$, (10.4) follows at once from (10.3). Assume (10.4) for a given $k$. Applying (2.6) and (2.7), we obtain

$$c(f_k(z_k)) = \gamma^N$$

and

$$d(f_k(z_k)) = d_N.$$

Hence, by (6.2),

$$s(f_k(z_k)) = (u(\gamma^N), d_N). \tag{10.5}$$

Now by (6.4) and the linearity of $f_k^i$ we have

$$f_{k+1}^i(z_{k+1}) = \sum_{y \in J_1 \times \ldots \times J_k}^{*} z_k(y) f_{k+1}^i(y, s(f_k(z_k)))$$

$$= \sum_{y \in J_1 \times \ldots \times J_k}^{*} z_k(y) f_{k+1}^i(y, (u(\gamma^N), d_N)) \quad \text{(by (10.5))}$$

$$= \overset{*}{\underset{y \in J_1 \times ... \times J_k}{\sum}} z_k(y) f^i_{k+1}(y^i, (u^i(\gamma^N), N)) \quad \text{(by (2.4) and (2.5))}$$

$$= \overset{*}{\underset{y \in J_1 \times ... \times J_k}{\sum}} z_k(y) \gamma^N \qquad\qquad \text{(by (10.3))}$$

$$= \gamma^N \underset{y \in J_1 \times ... \times J_k}{\sum} z_k(y)$$

$$= \gamma^N. \qquad\qquad\qquad \text{(by (2.2))}$$

This completes the induction and the proof of the lemma.

LEMMA 10.6    For $k, j \geqslant 0$ and $k \neq j$, $c(f_k(z_k))$ is statistically independent of $c(f_j(z_j))$.

*Proof*   By (10.3), the choice of strategies after each play depends only on the coalitions previously chosen. Thus the sequence of partitions of $N$ into coalitions is strictly determined. For a given play, the choice of $c$-strategy $B$-vectors depends only on this sequence (which is, in fact, a sequence of constant partitions), and *not* on the strategies previously chosen. This completes the proof.

LEMMA 10.7   $f$ is summable, and $H(f) = h$.

*Proof*   Applying (2.7) to (10.4), we obtain

$$c(f_k(z_k)) = \gamma^N,$$

whence by (6.6)

$$E_k = E(c(f_k(z_k))) = E(\gamma^N). \tag{10.8}$$

By Lemma 10.6, the strong law of large numbers[7] applies to each component of $c(f_k(z_k))$, hence to $c(f_k(z_k))$ itself, and hence also to $E_k$. By (10.8) and (6.7), we have

$$H_k = H(\gamma^N)$$

$$= h. \qquad \text{(by (10.1))}$$

Hence

$$H(f) = \lim_{k \to \infty} \frac{1}{k} \sum_{r=1}^{k} H_r \quad \text{(by (6.8))}$$

$$= \lim_{k \to \infty} \frac{1}{k} kh$$

---

7. See [2], p. 207, *The Kolmogoroff Criterion*, or p. 208, the Theorem.

$$= \lim_{k \to \infty} h$$
$$= h.$$

By Definition 6.9, the proof is complete.

LEMMA 10.9    If $v_k | R$ is a partition and $k > 0$, then

$$j \in e(f_k^i(v_1, \ldots, v_k)) \Longrightarrow f_k^j(v_1, \ldots, v_k) = f_k^i(v_1, \ldots, v_k).$$

*Proof*    From (10.3) and the left side of the conclusion, we deduce that $j \in v_k^i | R^i$. Hence since $v_k | R$ is a partition, it follows that

$$v_k^i | R^i = v_k^j | R^j,$$

and then the right side of the conclusion follows at once from (10.3).

In order to show that $f$ is a strong equilibrium $c$-point, we will first assume that it is not.

ASSUMPTION 10.10    There is a $B \subset N$ (which we will call $B*$) for which there is a supergame $c$-strategy vector $g$ for which

$$g^{N-B*} = f^{N-B*} \tag{10.11}$$

and there is a $B$*-vector $\varepsilon^{B*} > 0$ for which

$$\lim_{k \to \infty} \mathrm{Prob}_g(S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \text{ for some } r \geqslant k) > 0. \tag{10.12}$$

It will be our task to show that (10.10) leads to an absurdity.

Let $v = (v_1, \ldots, v_k, \ldots)$ be an arbitrary sequence for which $v_k \in J_k$, $k \geqslant 1$. Let $B(v)$ denote the set of $i$ for which there is a $k_i \geqslant 0$ such that

$$g_{k_i}^i(v_1, \ldots, v_{k_i}) \neq f_{k_i}^i(v_1, \ldots, v_{k_i}). \tag{10.13}$$

It follows from (10.11) and (10.13) that

$$B(v) \subset B*. \tag{10.14}$$

DEFINITION 10.18    For each $v$ and for $k \geqslant 0$, let $B_k(v)$ denote the set of $i$ for which there is a $k_i$ such that

$$0 \leqslant k_i \leqslant k$$

and such that (10.13) holds. Define

$$B_{-1}(v) = \varnothing. \tag{10.19}$$

For each $v$, there is clearly a $k(v)$ such that

$$B_k(v) = B(v), \quad k \geqslant k(v). \tag{10.20}$$

Furthermore, it is not difficult to see that (10.21): For each $v$, $B_k(v)$ is monotone increasing with $k$. We will in the sequel restrict use of the symbol $k(v)$ to the *first* $k(v)$ satisfying (10.20).

DEFINITION 10.22   We say that $v$ occurs with *positive probability* if for each $k > 0$,

$$z_k(g)(v_1, \ldots, v_k) > 0;$$

that is, if for each $k > 0$, and under the assumption that the supergame $c$-strategy vector $g$ is played, $(v_1, \ldots, v_k)$ occurs with positive probability.

LEMMA 10.23   If $v$ occurs with positive probability, then for each $k > 0$, $v_k | R$ is a partition.

*Proof*   By (6.4) and (10.22), $(v_1, \ldots, v_k)$ occurs with positive coefficient in

$$(v_1, \ldots, v_{k-1}, s(g_{k-1}(v_1, \ldots, v_{k-1}))).$$

Hence $v_k$ occurs with positive coefficient in $s(g_{k-1}(v_1, \ldots, v_{k-1}))$. Applying (6.2), we conclude that $v_k | R$ occurs with positive coefficient in $d(g_{k-1}(v_1, \ldots, v_{k-1}))$. Since by definition, $d(g_{k-1}(v_1, \ldots, v_{k-1}))$ is a probability combination of a number of partitions of $N$, it follows that $v_k | R$ must be a partition, which completes the proof.

LEMMA 10.24   Let $v$ occur with positive probability, and assume that (10.10) holds. For $k \geqslant 0$ and $i \in N - B_k(v)$, we have

$$g_k^i(v_1, \ldots, v_k) = \gamma^{N - B_{k-1}(v)}.$$

*Proof*   We will write $B$ instead of $B(v)$, and for each $k$, $B_k$ instead of $B_k(v)$. The proof is by induction on $k$. First let $k = 0$. If $i \in N - B_0$, then it follows from Definition 10.18 that there is no $k_i$ satisfying $k_i = 0$ and (10.13). In other words, for $i \in N - B_0$, we have

$$g_0^i = f_0^i.$$

Hence by (10.3), we have for $i \in N - B_0$,

$$g_0^i = \gamma^N$$
$$\quad = \gamma^{N - B_{-1}} \quad \text{(by (10.19))},$$

which establishes (10.24) for $k = 0$.

Next, suppose we have established (10.24) for some $k \geqslant 0$ and all preceding $k$. For $i \in B_k - B_{k-1} = (N - B_{k-1}) \cap B_k$, we must have, by Defi-

nition 10.18, that (10.13) holds for $k_i = k$. Hence

$$g_k^i(v_1, \ldots, v_k) \neq f_k^i(v_1, \ldots, v_k), \quad i \in B_k - B_{k-1}. \tag{10.25}$$

Now if $k = 0$, it follows from (10.19), (10.25), and (10.3) that

$$g_0^i \neq \gamma^{N - B_{-1}}, \quad i \in B_0 - B_{-1}. \tag{10.26}$$

If $k > 0$, we have for $j \in B_k - B_{k-1}$ that

$$j \in N - B_{k-1} = e(g_k^i(v_1, \ldots, v_k)), \quad i \in N - B_k \quad \text{(by induction hypothesis)}$$
$$= e(f_k^i(v_1, \ldots, v_k)), \quad i \in N - B_k \quad \text{(by (10.18))}.$$

Applying Lemma 10.9 and Lemma 10.23, we obtain

$$f_k^i(v_1, \ldots, v_k) = f_k^j(v_1, \ldots, v_k), \quad k > 0, \quad i \in N - B_k, \quad j \in B_k - B_{k-1}. \tag{10.27}$$

From (10.18) and the induction hypothesis, we obtain

$$f_k^i(v_1, \ldots, v_k) = g_k^i(v_1, \ldots, v_k) = \gamma^{N - B_{k-1}}, \quad i \in N - B_k, \quad k > 0.$$

Combining this with (10.27), we obtain

$$f_k^j(v_1, \ldots, v_k) = \gamma^{N - B_{k-1}}, \quad j \in B_k - B_{k-1}, \quad k > 0,$$

i.e.,

$$f_k^i(v_1, \ldots, v_k) = \gamma^{N - B_{k-1}}, \quad i \in B_k - B_{k-1}, \quad k > 0.$$

From this and (10.25) we deduce, for $k > 0$, that

$$g_k^i(v_1, \ldots, v_k) \neq \gamma^{N - B_{k-1}}, \quad i \in B_k - B_{k-1}; \tag{10.28}$$

the same result for $k = 0$ is given by (10.26). Hence (10.28) holds for $k \geqslant 0$.

We certainly have

$$g_k^i(v_1, \ldots, v_k) \neq \gamma^{N - B_{k-1}}, \quad i \in B_{k-1},$$

for otherwise $i \notin e(g_k^i(v_1, \ldots, v_k))$, contradicting the membership of $g_k^i$ in $T^i$. Combining this with the induction hypothesis and (10.28), we obtain

$$g_k^i(v_1, \ldots, v_k) \begin{cases} = \gamma^{N - B_{k-1}}, & i \in N - B_k \\ \neq \gamma^{N - B_{k-1}}, & i \in B_k. \end{cases}$$

It follows that

$$d^i(g_k(v_1, \ldots, v_k)) = N - B_k, \quad i \in N - B_k,$$

and in particular, since $N - B_{k+1} \subset N - B_k$, we have

$$d^i(g_k(v_1, \ldots, v_k)) = N - B_k, \quad i \in N - B_{k+1}. \tag{10.29}$$

Now for $i \in N - B_{k+1}$, we have by Definition 10.18 that

$$g^i_{k+1}(v_1, \ldots, v_{k+1}) = f^i_{k+1}(v_1, \ldots, v_{k+1}) = \gamma^{v^i_{k+1}|R^i}. \tag{10.30}$$

By reasoning similar to that used in the proof of (10.23), we obtain that $v^i_{k+1}|R^i$ occurs with positive coefficient in $d^i(g_k(v_1, \ldots, v_k))$. Combining this with (10.29) and (10.30), we obtain that for $i \in N - B_{k+1}$

$$g^i_{k+1}(v_1, \ldots, v_{k+1}) = \gamma^{N-B_k}.$$

This completes the induction and establishes (10.24).

From (10.24) and (10.20) we deduce that if $v$ occurs with positive probability, then

$$g^i_k(v_1, \ldots, v_k) = \gamma^{N-B(v)}, \quad i \in N - B(v), \quad k > k(v),$$

whence

$$c(g_k(v_1, \ldots, v_k)) = (\gamma^{N-B(v)}, \ c^{B(v)}(g_k(v_1, \ldots, v_k))), \quad k > k(v).$$

Setting

$$c^B(g_k(v_1, \ldots, v_k)) = \gamma^B_k, \quad B \subset N, \quad k > k(v),$$

we obtain

$$c(g_k(v_1, \ldots, v_k)) = (\gamma^{B(v)}_k, \ \gamma^{N-B(v)}), \quad k > k(v). \tag{10.31}$$

LEMMA 10.32   For each $m \geqslant 0$, there is a positive integer $k$, such that for all $v$, $r$, and $h_1, \ldots, h_m$ for which

$$r \geqslant k \tag{10.33}$$

and

$$h_j \in H(P), \quad j = 1, \ldots, m, \tag{10.34}$$

and

$$S^{B*}_r(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \tag{10.35}$$

we have

$$r > m \tag{10.36}$$

and

$$\frac{1}{r}\left(\sum_{j=1}^{m} h_j^{B*} + \sum_{j=m+1}^{r} H_j^{B*}(v)\right) \geqslant H^{B*}(f) + \frac{1}{3}\varepsilon^{B*}. \tag{10.37}$$

*Proof*   First of all, we may choose

$$k > m,$$

so that (10.36) is satisfied. Next, choose

$$k \geqslant \max_{i \in B*} \frac{3m}{\varepsilon^i} \max_{p \in P} H^i(p). \tag{10.38}$$

We then have

$$\left| S_r^{B*}(v) - \frac{1}{r}\left(\sum_{j=1}^{m} h_j^{B*} + \sum_{j=m+1}^{r} H_j^{B*}(v)\right)\right|$$

$$= \left|\frac{1}{r}\left(\sum_{j=1}^{r} H_j^{B*}(v)\right) - \frac{1}{r}\left(\sum_{j=1}^{m} h_j^{B*} + \sum_{j=m+1}^{r} H_j^{B*}(v)\right)\right| \qquad \text{(by (6.11))}$$

$$= \frac{1}{r}\left|\sum_{j=1}^{m}(H_j^{B*}(v) - h_j^{B*})\right|$$

$$\leqslant \frac{1}{k}\sum_{j=1}^{m}\left|H_j^{B*}(v) - h_j^{B*}\right| \qquad \text{(by (10.33))}$$

$$\leqslant \frac{1}{k}\sum_{j=1}^{m} 2\max_{p \in P} H^{B*}(p) \qquad \text{(by (10.34))}$$

$$= \frac{2m}{k}\max_{p \in P} H^{B*}(p)$$

$$\leqslant \frac{2}{3}\varepsilon^{B*} \qquad \text{(by (10.38))}$$

Combining this with (10.35), we obtain (10.37). This completes the proof.

COROLLARY 10.39   For each $m \geqslant 0$, the following holds:
  For all sufficiently large $k$, we have that for all $v$, $r$, and $h_1, \ldots, h_m$ for which (10.33), (10.34), and (10.35) hold, (10.36) and (10.37) follow.

*Proof*   Because if (10.32) is true for a given $k$, it is certainly true for all larger $k$.

LEMMA 10.40   For each $m \geqslant 0$ and $B \subset N$, including the null set, we have

$$\lim_{k \to \infty} \text{Prob}_g((S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \text{ for some } r \geqslant k)$$

and $k(v) = m$ and $B(v) = B) = 0$.

*Proof*   First let $B = \varnothing$. Then for the $v$ we are considering (inside $\mathrm{Prob}_g$), we have

$$B(v) = \varnothing. \tag{10.41}$$

Since $B_k(v)$ is monotone increasing with $k$ and always contained in $B(v)$, it follows from (10.41) that for all $k$,

$$B_k(v) = \varnothing.$$

Applying (10.24), we obtain for $k \geqslant 0$ and $i \in N$,

$$g_k^i(v_1, \ldots, v_k) = \gamma^N. \tag{10.42}$$

Let $\theta$ be a supergame $c$-strategy vector given for all $v$ by

$$\theta_k^i(v_1, \ldots, v_k) = \gamma^N, \quad k \geqslant 0, \quad i \in N. \tag{10.43}$$

From (10.42) and (10.43), we obtain

$$\lim_{k \to \infty} \mathrm{Prob}_g((S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \qquad \text{for some } r \geqslant k)$$
$$\text{and } k(v) = m$$
$$\text{and } B(v) = B).$$

$$= \lim_{k \to \infty} \mathrm{Prob}_\theta((S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \quad \text{for some } r \geqslant k)$$
$$\text{and } k(v) = m$$
$$\text{and } B(v) = B)$$

$$\leqslant \lim_{k \to \infty} \mathrm{Prob}_\theta(S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \text{ for some } r \geqslant k)$$

$$\leqslant \lim_{k \to \infty} \mathrm{Prob}_\theta(S_r^i(v) \geqslant H^i(f) + \varepsilon^i \text{ for some } r \geqslant k)$$

$$= \lim_{k \to \infty} (1 - \mathrm{Prob}_\theta(S_r^i(v) < H^i(f) + \varepsilon^i \text{ for all } r \geqslant k))$$

$$= 1 - \lim_{k \to \infty} \mathrm{Prob}_\theta(S_r^i(v) < H^i(f) + \varepsilon^i \text{ for all } r \geqslant k)$$

$$\leqslant 1 - \lim_{k \to \infty} \mathrm{Prob}_\theta(|S_r^i(v) - H^i(f)| < \varepsilon^i \text{ for all } r \geqslant k), \tag{10.44}$$

where $i$ is an arbitrary member of $B*$.

Now when $v$ is distributed according to $\theta$, we deduce from (10.43) and (6.10) that the $H_j^i(v)$ are independent random variables whose mean is $H^i(\gamma^N)$. Applying (6.11), (10.1), (10.7), and the strong law of large numbers, we obtain that the second term of the right side of (10.44) is 1,

whence we deduce that the left side vanishes. This completes the proof in the case $B = \varnothing$.

Next, suppose $B \neq \varnothing$. From (10.14) we obtain that $B(v) \subset B*$. Thus if $B \not\subset B*$, then $\mathrm{Prob}_g(B(v) = B) = 0$, and our lemma is already proved. Thus we may assume without loss of generality that

$$B \subset B*. \tag{10.45}$$

For $j = 1, \ldots, m$, suppose the random variable $w_j^B$ to be distributed according to an arbitrary but fixed member of $C^B(J_j^B)$, which we will call $\gamma_j^B$, while $w_j^{N-B}$ is distributed according to $\gamma^{N-B}$. Set

$$h_j = H(w_j), \quad j = 1, \ldots, m. \tag{10.46}$$

Applying Corollary 10.39, we obtain

$$\lim_{k \to \infty} \mathrm{Prob}_g((S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \quad \text{for some } r \geqslant k)$$
$$\text{and } k(v) = m$$
$$\text{and } B(v) = B)$$

$$\leqslant \lim_{k \to \infty} \mathrm{Prob}\left(\left(\frac{1}{r}\left(\sum_{j=1}^m h_j^{B*} + \sum_{j=m+1}^r H_j^{B*}(v)\right) \geqslant H^{B*}(f) + \frac{1}{3}\,\varepsilon^{B*}\right.$$
$$\text{for some } r \geqslant k)$$
$$\text{and } k(v) = m$$
$$\text{and } B(v) = B), \tag{10.47}$$

where $v_{m+1}, \ldots, v_r$ are distributed according to $x_{m+1}(g), \ldots, x_r(g)$, while $h_j$ is given by (10.46). If for $j = k+1, \ldots, r$, we will let $w_j$ be a random variable whose distribution is $(\gamma_j^B, \gamma^{N-B})$, where the $w_j^B$ may be dependent on $w_\ell$ with $\ell < j$ but the $w_j^{N-B}$ must be independent of all $w_\ell$ with $\ell \neq j$, then we may conclude from (10.31) and (6.10) that the right side of (10.47) is

$$\leqslant \lim_{k \to \infty} \mathrm{Prob}\left(\frac{1}{r}\sum_{j=1}^r H^{B*}(w_j) \geqslant H^{B*}(f) + \frac{1}{3}\,\varepsilon^{B*} \text{ for some } r \geqslant k\right)$$

$$\leqslant \lim_{k \to \infty} \mathrm{Prob}\left(\frac{1}{r}\sum_{j=1}^r H^B(w_j) \geqslant H^B(f) + \frac{1}{3}\,\varepsilon^B \text{ for some } r \geqslant k\right)$$
$$\tag{10.48}$$

by (10.45). Now there is certainly a supergame $c$-strategy vector $\beta$ for which

$$c(\beta_{j-1}(z_{j-1}(\beta))) = (\gamma_j^B, \gamma^{N-B}), \quad j \geqslant 1,$$

whence the right side of (10.48)

$$= \lim_{k \to \infty} \mathrm{Prob}_\beta(S_r^B(v) \geqslant H^B(f) + \tfrac{1}{3}\varepsilon^B \text{ for some } r \geqslant k). \tag{10.49}$$

Applying (10.2), (10.7) and Lemma 8.3, we obtain that the right side of (10.49) vanishes. Combining this with (10.48) and (10.47), we obtain the lemma.

LEMMA 10.50

$$\lim_{k \to \infty} \mathrm{Prob}_g(S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \text{ for some } r \geqslant k) = 0.$$

*Proof*   Letting $\bigcup$ stand for mutually exclusive disjunction, we have

$$1 = \mathrm{Prob}_g\left( \bigcup_{m=o}^{\infty} k(v) = m \right)$$

$$= \sum_{m=o}^{\infty} \mathrm{Prob}_g(k(v) = m)$$

$$= \sum_{m=o}^{\infty} \mathrm{Prob}_g\left( k(v) = m \text{ and } \bigcup_{B \subset N} (B(v) = B) \right)$$

$$= \sum_{m=o}^{\infty} \sum_{B \subset N} \mathrm{Prob}_g(k(v) = m \text{ and } B(v) = B). \tag{10.51}$$

Let $Q(v, k)$ stand for the expression

$$S_r^{B*}(v) \geqslant H^{B*}(f) + \varepsilon^{B*} \text{ for some } r \geqslant k.$$

Then

$$\lim_{k \to \infty} \mathrm{Prob}_g(Q(v, k))$$

$$= \lim_{k \to \infty} \mathrm{Prob}_g(Q(v, k) \text{ and } \bigcup_{m=o}^{\infty} (k(v) = m) \text{ and } \bigcup_{B \subset N} (B(v) = B))$$

$$= \lim_{k \to \infty} \sum_{m=o}^{\infty} \sum_{B \subset N} \mathrm{Prob}_g(Q(v, k) \text{ and } k(v) = m \text{ and } B(v) = B)$$

$$= \sum_{m=o}^{\infty} \sum_{B \subset N} \lim_{k \to \infty} \mathrm{Prob}_g(Q(v, k) \text{ and } k(v) = m \text{ and } B(v) = B)$$

(since the series is dominated by the series on the right side of (10.51))

$$= \sum_{m=o}^{\infty} \sum_{B \subset N} 0 \quad \text{(by Lemma 10.40)}$$

$$= 0.$$

This completes the proof of Lemma 10.50.

Lemma 10.50 contradicts (10.12) and thus establishes the truth of Theorem 3.

COROLLARY 4   $H(A_c) = H(S_c) = H(\tilde{S}_c)$.

*Proof*   Combine (7.4), Theorem 2 and Theorem 3.

---

## 11   Existence of Acceptable Points

All two-person games, zero-sum or not, have acceptable points. The proof, which is not difficult, will be given in a subsequent paper of this series devoted exclusively to the application of this theory to two-person games. Unfortunately, though, when we go beyond two-person games, we find games in which there are no acceptable points. It is instructive to examine an example of such a game.

$G$ is a three-person, zero-sum game. The three persons form a community and at regular intervals hold elections for mayor and vice-mayor. The mayor draws a salary of 2; the vice-mayor draws a salary of 1; the remaining player pays the salaries.

Let $h$ be an arbitrary payoff vector. Since $G$ is zero-sum, not all the players can have a positive payoff. Without loss of generality, let

$$h^1 \leqslant 0. \tag{11.1}$$

Now the minimum payoff to 1, even in pure strategies, is $-3$; hence since $G$ is zero-sum, it follows that

$$h^2 + h^3 \leqslant 3. \tag{11.2}$$

From (11.2) it follows that at least one of 2 and 3 gets at most $1\frac{1}{2}$; so we may write without loss of generality

$$h^2 \leqslant 1\tfrac{1}{2}. \tag{11.3}$$

Now by agreeing to vote for 2 for mayor and for 1 for vice-mayor, 1 and 2 can both increase their payoffs; and there is nothing that 3 can do to prevent this.

Actually, the game $G$ just described does not *deserve* an acceptable point. In a long sequence of plays of $G$, a steady state will never be reached; no sooner do we start settling down to one, than two of the players will see that it is relatively disadvantageous for them, and will certainly not be willing to agree to it for the remainder of the superplay. The game is *inherently unstable*. The same argument holds for all games not possessing acceptable points. Thus in seeking a steady state for the supergame, it would seem that we may restrict our-
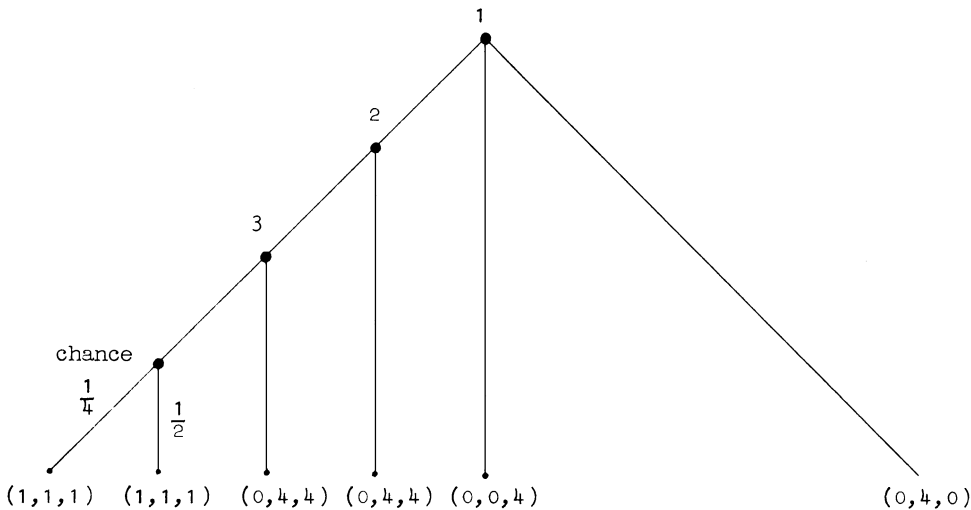
**Figure 1**
The payoff vector at each terminal gives the payoff to players 1, 2, and 3 in that order.

selves to games possessing acceptable points. These are called *stable* games.

---

## 12   The Expected Payoff

In order to see why the notion defined by (6.8) is essentially misleading, let us consider the three-person game $G$ of perfect information which is given in extensive form in Figure 1.

The payoff vector at each terminal gives the payoff to players 1, 2, and 3 in that order.

Player 1 has three strategies, $\ell^1$, $c^1$, and $r^1$ (left, center, and right). The other players have two strategies each; they will be denoted by $\ell^j$ and $r^j$, where $j = 0, 2$, and 3. Chance[8] is denoted by 0.

The point $(\ell^1, \ell^2, \ell^3)$, yielding a payoff of $(1, 1, 1)$, is acceptable. No player or group of players can obtain a higher payoff. However, if 1 plays the supergame $c$-strategy given by (10.3), there is a supergame $c$-strategy (2,3)-vector that will yield an *expected* payoff of $(0, 2, 2)$ rather than

---

8. It may be objected that we have excluded chance from consideration in this paper. Actually, chance is used in this example only to provide, as conveniently and cheaply as possible, a random device on which players 2 and 3 will be able to peg their choices. We could eliminate chance at the cost of complicating the example considerably. In order to exhibit the principle involved as clearly as possible, we have retained chance; but it is important to remember that it does not form an essential part of the example.

(1, 1, 1). This supergame $c$-strategy (2,3)-vector, which we will call the "plan," may be described as follows: On the first play, both players make the "orthodox" choice, i.e., they concur with player 1 in choosing $(\ell^1, \ell^2, \ell^3)$. On all subsequent plays, we distinguish two cases:

1. If chance played $\ell^o$ on the first play, then player 2 makes the "orthodox" choice, i.e., that given by (10.3), while player 3 plays $r^3$.

2. If chance played $r^o$ on the first play, then player 3 makes the "orthodox" choice, while player 2 plays $r^2$.

According to (10.3), player 1 will respond to this treachery with a wounded roar and proceed to wreak dire vengeance. In case 1 he will finish off player 3 by playing $r^1$; in case 2 he will play $c^1$, thus avenging himself on the treacherous player 2. In each case, he is convinced that the orthodox player has been loyal all along, and therefore will not feel too unhappy about his getting 4 while player 1 himself gets nothing. The expected payoff is (0, 2, 2), an improvement for both 2 and 3. For 2 and 3, this sounds like a good set-up, at least on paper. Who wouldn't want an expected payoff of 2 instead of 1? Let's see how it works out in practice.

Imagine that chance has played $\ell^o$ on the first play. For the second play, player 1, all unsuspecting, suggests $(\ell^1, \ell^2, \ell^3)$, and player 2 concurs, according to plan. Player 3 is now faced with the choice of going along with the plan and playing $r^3$, thus bringing down upon himself the eternal wrath of player 1, or of concurring in the choice of $(\ell^1, \ell^2, \ell^3)$, thus possibly making player 2 angry at him, but still assuring himself of a steady income of 1. It won't take him long to make his choice. He has nothing to gain and everything to lose by going along with the plan; he will get nothing but 0 for the rest of eternity for his pains, while player 2 enjoys a steady payoff of 4. Player 3 will not now be consoled by the fact that if chance had played the other way, *he* would have been getting 4 while player 2 would have been getting 0; in no sense does this constitute an incentive for him to follow the plan as things stand now. The point is that in a single play of a game, it may be worthwhile for each of the parties to plan to co-operate, even if it hurts, because in future plays, *he* may be the beneficiary. But in a superplay, there *are* no future superplays; *this is it*; by going along with the plan, player 3 will be ruining himself not just for one play, but for always. It just isn't worthwhile, and no player would do it. Similar remarks apply to player 2. The plan does yield a higher expected payoff, but will never be carried out.

It is possible to define strong equilibrium $c$-points using an appropriate formula involving expected payoffs in place of 7.2. If this is done, the first half of the main theorem, the analogue of Theorem 2, remains true.

However, the second half, the analogue of Theorem 3, fails. The game $G$ that we have been discussing provides the counterexample. If we start out with $h = (1, 1, 1)$, then not only does the $f$ given by (10.3) fail, as we have shown, but it can be shown that there is then no strong equilibrium $c$-point (in the sense of expected payoffs) whose payoff is $h$.

## 13   Enforceability of Agreements

There are two kinds of agreements involved in a cooperative supergame. One is the formal agreement that is contracted prior to a given play by the members of a coalition $B$, who agree to play a certain correlated strategy $B$-vector on that play. This agreement expires immediately after the play in question is completed; its validity does not extend to any subsequent plays. The other is the informal "gentlemen's agreement" in which the players undertake informally to make certain choices or join certain coalitions on each of a number of plays, or even for the entire superplay; they do not form a part of the formal mathematical structure introduced in Section 6. These are the agreements that are involved in the supergame strategy vectors that we have used in our proofs. When in heuristic discussions above we referred to "orthodoxy" or "deviationism," it was the second kind of agreement to which we were referring. To distinguish between the two kinds of agreements, we will call the first formal, the second informal.

As mentioned in Section 3, together with the formal agreement there comes a formal enforcement mechanism. Put completely precisely, if for the $k^{\text{th}}$ play the players choose the vector $t \in T$, then each $i \in N$ *must* become a member of $d^i(t)$, and the coalition $d^i(t)$ *must* then play $(t^i)^{d^i(t)}$. If the reader wishes, he may consider that once the choice of $t$ has been made, the remainder of the play is no longer in the hands of the individual players, but in the hands of an umpire. Such an extreme interpretation is unnecessary for practical purposes; we refer to it only to make clear the intuitive meaning of formal enforcement.

On the other hand, the informal agreement can be enforced in no such formal way. Violation of such an agreement can be "prevented" only by the threat or implied threat of retaliation on subsequent plays. The exact nature of this "pseudo-enforceability" is of course the crux of what we have been investigating in the foregoing sections.

In the light of these remarks, it is an interesting fact that if in the supergame we abandon formal enforceability for the formal agreements, the main theorem remains true. The possibility of retaliation is sufficient to enforce even the formal agreements, without any necessity for formal

enforcement apparatus. We have not proved this statement in the foregoing; the structure introduced in Section 6 was constructed with formal enforceability of formal agreements in mind. A structure permitting the violation of formal agreements would be somewhat more complicated.[9] However, upon a little reflection, the reader will be able to convince himself of the truth of the statement, at least intuitively.

Of course, when considering just a single play rather than the supergame, as in the definition of acceptability, then there are only formal agreements and these are formally enforceable.

---

## References

1. Nash, J. F., "Non-cooperative games," Annals of Mathematics 54 (1951), pp. 286–295.

2. Feller, W., An Introduction to Probability Theory and its Applications, John Wiley and Sons, 1950.

3. Kuhn, H. W., "Extensive games and the problems of information," Annals of Mathematics Study No. 28 (Princeton, 1953), pp. 193–216.

4. Blackwell, D., "An analog of the minimax theorem for vector payoffs," Pacific J. Math. 6 (1956), pp. 1–8.

5. Thompson, F. B., "Equivalence of games in extensive form," The RAND Corporation, Research Memorandum 759, 1952.

6. Luce, R. D., and Raiffa, H., Games and Decisions, John Wiley and Sons, 1957.

7. Von Neumann, J., and Morgenstern, O., Theory of Games and Economic Behavior, Princeton, 1944, 3rd ed., 1953.

8. Raiffa, H., "Arbitration schemes for generalized two-person games," Annals of Mathematics Study No. 28 (Princeton, 1953), pp. 361–387.

9. Roughly speaking, each play would be described by a choice of a $t \in T$ as in Section 6, followed by a chance choice of a pure strategy vector in accordance with $c(t)$, followed by some move permitting the players to renege on their commitments if they wish.