

Discussion Paper No. 786

INFORMATION LEAKAGE FORCES COOPERATION

by

Akihiko Matsui

Northwestern University

May 1988

The author gratefully acknowledges many useful comments from Itzhak Gilboa, Hisao Hisamoto, Ehud Kalai, and Dov Samet. The author has also enjoyed conversations with Mamoru Kaneko and Mikio Nakayama. All errors, of course, are my own.

Correspondence: Akihiko Matsui
Managerial Economics and Decision Sciences
J. L. Kellogg Graduate School of Management
Northwestern University
Evanston, IL 60208
U.S.A.

Abstract

This paper considers a two-person repeated game in which there is a small probability of espionage, i.e., that one or both of the players will be informed of the other's supergame strategy and have a chance to revise his strategy based on this information before the game begins. It is shown that in such a game any subgame perfect equilibrium pair of payoffs is Pareto efficient provided that the probability of espionage is small enough.

In view of the "Folk Theorem," several attempts have been made to shrink the set of outcomes that are supported by equilibria in an infinitely repeated game. Two main approaches are to create new solution concepts which allow "collective deviation," and to introduce bounded rationality considerations. The present paper assumes neither "collective deviation" nor bounded rationality to derive the result.

1. INTRODUCTION

Since the seminal paper by Aumann [3], repeated games have drawn attention of many theorists.^{1/} The underlying observation that has motivated those studies is that in a one-shot game an undesirable outcome is often supported by an equilibrium, and moreover, it may happen that only inefficient outcomes are equilibrium outcomes as in the case of classical prisoners' dilemma example. This unfavorable situation is partially reconciled when one considers the repetition of the one-shot game. In a repeated game players can make their actions depend upon what they observed. This implies that players can punish those who took unfavorable action in some preceding periods. In this situation players may be reluctant to deviate from a socially favorable equilibrium for fear of punishment in later periods in spite of an instant extra benefit the deviation may yield. Because of this mechanism, the set of equilibrium outcomes may differ from the simple repetition of the equilibrium outcomes of the one-shot game. The repeated games are classified into two categories according to the length of repetition of a game: finitely repeated games and infinitely repeated games. The results for these two types of repeated games differ.^{2/} As far as infinitely repeated games are concerned, the most important property is the Folk theorem, which basically states that every feasible pair of individually rational payoffs can be attained by a noncooperative equilibrium.^{3/}

In those papers which derive the Folk theorem, Nash equilibrium or its refined concepts such as subgame perfect equilibrium have been used as the criteria for equilibria (see e.g. Rubinstein [11], Abreu [1]). The standard

assumption underlying these results is that each player can observe only the actions taken by the other players, rather their whole supergame strategies, and hence a player cannot condition his choices on these strategies. The purpose of this paper is to analyze what will occur if there is a small chance of leakage of information on players' actual supergame strategies. To this aim, I shall build a two-person repeated game in which there is a small probability of espionage, i.e., that one or both of the players will be informed of the other's supergame strategy and have a chance to revise his strategy based on this information before the game begins. It is shown that in such a game any subgame perfect equilibrium pair of payoffs is Pareto efficient for sufficiently small probability of espionage. Furthermore, in the case of one-sided espionage, this result is shown to hold for any positive probability of espionage.

The above question is worth being raised because in many industries spying activities, which may be either legal or illegal, have important roles, and by means of these activities there is always a small chance of information leakage. Suppose that a firm's strategy is determined in meetings of board of directors and filed in meeting protocols and other written statements. This strategy may leak to its opponent by means of industrial espionage with at least small probability. If the firm (board of directors) knows that this is the case, then it takes that possibility into account in determining its strategy. The results may be quite different from those without the possibility of information leakage.

From a theoretical point of view, attempts have been made to shrink the set of outcomes that are supported by equilibria. Two main streams can be discerned in this literature. One of them is to cope with the problem by

creating new solution concepts which allow "collective deviation." The earliest attempt was made by Aumann [3], though his main motivation was to build a non-cooperative foundation for cooperative game theory. Another effort is made by Pearce [10], who deals with the problem by using a new equilibrium concept called renegotiation-proof equilibrium.

Those papers cope with the situation where unlimited communication is available between participants of the game and assume that communication makes it possible for players to deviate simultaneously. As opposed to those papers, the present paper assumes only a small probability of information leakage without any negotiation procedure. Furthermore, this paper does not assume any possibility of "collective deviation." Still the result is that cooperation emerges.

It is worth noting here that this paper is related to Kalai [7] though the latter does not deal with repeated games. He built a formal model for preplay communication process and derived the result that only Pareto efficient outcomes are supported by a noncooperative (subgame perfect) equilibrium, which shares the basic direction of research with this paper.

The other stream is to deal with the problem by using finite automata. Rubinstein [12] restricts his attention to infinitely repeated prisoners' dilemma and explicitly introduces the cost of implementation of strategies. Under the assumption that each player has lexicographical preferences in which the repeated game payoff matters first, and the complexity cost measured by the number of internal states in the automaton matters second, it is shown that the Folk Theorem does not hold any more. Abreu and Rubinstein [2] extend the model to other cases and obtain the similar results. Here, however, the Nash equilibrium in a one-shot game still

remains an equilibrium.

Another type of analysis is made by Kreps, Milgrom, Roberts, and Wilson [9] and Aumann and Sorin [5].^{4/} They show that if there is a small probability that the opponent is irrational, only one Pareto efficient outcome appears as an equilibrium. The point of these papers is not only that the opponent is not rational with small probability, but also that this irrationality should be restricted to a certain class of automata. In Kreps et al. [9], the only possible irrationality permitted in the model is "tit for tat", which starts with cooperation and after that follows the previous action of the other player. In this sense cooperation is incorporated in the model. In fact, depending upon the choice of machine representing irrationality, every feasible pair of strictly individually rational payoffs emerges as an equilibrium in a long but finitely repeated game. On the other hand, Aumann and Sorin [5] sophisticate the model by assuming that a player faces all the finite recall automata with positive probability. Here, however, one of the "simplest" strategies, grim-trigger strategy, is not permitted. Moreover, their result is restricted to a specific class of games in which the Pareto efficient set is a singleton.

The driving force for cooperative outcomes in the present paper is totally different from that of those two papers. In this paper players are perfectly rational, but with small probability one of them has a chance to revise his strategy after the opponent's strategy has been revealed to him. In very bold strokes, the logic of the main result may be explained as follows: Consider for simplicity the case of one-sided espionage, say, where player 2 may be informed of player 1's strategy. Next, consider a pair of strategies which yield a Pareto dominated pair of payoffs. Player 1 can

deviate to a strategy which expects a certain "signal" from player 2. If the signal is not received, player 1 will stick to his original strategy. if it is received, however, player 1 will switch to a strategy which corresponds to a Pareto dominating pair of payoffs. Thus, if player 2 happens to know player 1's strategy, his best response is to give the required signal and switch to the better payoff. Note that at equilibrium player 2 will give the signal and cooperate not only in the revision node after observing player 1's strategy but also from the very beginning.

It is worth mentioning that in order to obtain this result one has to assume that there is a certain cost of revision of strategies. Introducing the revision cost, one has to specify the trade-off between this cost and the supergame payoff. In the basic model, I deal with lexicographical preference, but later on it will be shown that a similar result can be obtained for other preferences in which the revision cost is not infinitesimal.

The contents of this paper are as follows. In Section 2 I shall present the model of one-sided espionage and lexicographical preferences. Section 3 illustrates the logic of the main result. In Section 4 I shall define an equilibrium to this game and find the set of outcomes which can be supported by equilibria. In Section 5 the result of Section 4 will be extended to the cases of two-sided espionage and positive revision cost. Section 6 concludes the paper.

2. MODEL

Supergame

Let $G = \langle S_1, S_2, \pi_1, \pi_2 \rangle$ be a two person game in normal form where S_i is a finite set of actions for player i ($i=1,2$), and $\pi_i: S_1 \times S_2 \rightarrow \mathbb{R}$ is the payoff function for player i . We assume that the number of actions in each S_i , $|S_i|$ ($i=1,2$), is greater than one. An infinite sequence of G is called a supergame of G and denoted by

$$G^\infty = (G_t)_{t=1}^\infty$$

where $G_t = G$. The supergame of G , G^∞ , is of standard information if every player can observe at t the actions of the other player which are taken before t . We confine our attention to a supergame with standard information. A strategy of the i -th player in a supergame with standard information is a sequence of functions $(f_i^t)_{t=1}^\infty$ such that $f_i^1 \in S_i$ and $f_i^t: S^{t-1} \rightarrow S_i$ for $t > 1$ where $S = S_1 \times S_2$. We also confine our attention to pure strategies. We denote by \mathcal{F}_i the set of all the supergame strategies of player i ($i=1,2$). Given the strategies chosen by both players, the supergame is played as follows. At the first period, the pair of actions is (f_1^1, f_2^1) . At the second period, the pair of actions is $(f_1^2(f_1^1, f_2^1), f_2^2(f_1^1, f_2^1))$. In this manner, a sequence of pairs of actions in G $((s_1^t, s_2^t))_{t=1}^\infty$ is deterministically generated by the pair of supergame strategies. We denote by $s_i^t(f_1, f_2)$ the action in G of player i at t determined by (f_1, f_2) . The supergame payoff of player i defined on $\mathcal{F}_i \times \mathcal{F}_j$ is assumed to be basically the limit of means:

$$\Pi_i(f_i, f_j) = \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \pi_i(s_1^t(f_1, f_2), s_2^t(f_1, f_2)), \quad i=1,2.$$

For later use, we introduce the following notations. Let $s_i \in S_i$ be such that

$$s_i \in \operatorname{argmin}_{s_i \in S_i} \max_{s_j \in S_j} \pi_j(s_1, s_2).$$

Let $\pi_1 = \max_{s_1 \in S_1} \pi_1(s_1, s_2)$ be the minimum individually rational payoff for player 1. Player 2 can impose the payoff π_1 for player 1 no matter what

strategy player 1 takes. The payoff π_2 is defined in a similar way. A strictly individually rational payoff for player 1 is a number which is greater than π_1 .

Information leakage game

Next, we describe a information leakage game (leakage game for short) in which two players choose supergame strategies. The leakage game consists of three stages. In the first stage, both players determine their supergame strategies simultaneously. In the second stage, Nature chooses one of the two alternatives with probability $1-\epsilon$ and the other with probability ϵ where we assume that ϵ is strictly between zero and unity. The first alternative (that with probability $1-\epsilon$) results in the supergame played by the supergame strategies which both players chose in the first stage. The second alternative brings player 2 to the third stage. In the third stage, player 2 is informed of player 1's strategy and may revise his supergame strategy. This revision is not restricted by what player 2 chose in the first stage. Player 1 is assumed to keep the strategy that he chose in the first stage. After the third stage ends, the supergame by the strategies determined by player 1 and 2 starts.

Strategies of player 1 and player 2 in the leakage game are described by a triple $(f_1, f_2; F)$ where $f_1 \in \mathcal{F}_1$, $f_2 \in \mathcal{F}_2$, and $F: \mathcal{F}_1 \times \mathcal{F}_2 \rightarrow \mathcal{F}_2$. A pair of supergame strategies (f_1, f_2) is chosen in the first stage, and a typical value of F , $F(f_1, f_2)$, is a supergame strategy taken by player 2 if f_1 and f_2 are chosen by player 1 and player 2 respectively, and if the revision node for player 2 is reached. The strategies in the leakage game may be called meta strategies in the sense that they determine supergame strategies. We

will refer to triples of supergame strategies $(f_1, f_2; f_2'')$ $\in \mathcal{F}_1 \times \mathcal{F}_2^2$ as outcomes of the leakage game. $(f_1, f_2; F(f_1, f_2))$ is an outcome induced by strategies of the leakage game $(f_1, f_2; F)$.

Preference relations of the players are defined on the set of outcomes of the leakage game. Player 1's preference relation is simply the ordering of his expected supergame payoff:

$$(1-\epsilon)\Pi_1(f_1, f_2) + \epsilon\Pi_1(f_1, F(f_1, f_2)).$$

We assume that the player 2's preference relation is lexicographical in the following sense. Player 2 prefers $(f_1, f_2; f_2'')$ to $(f_1', f_2'; f_2''')$ if

$$(1-\epsilon)\Pi_2(f_2, f_1) + \epsilon\Pi_2(f_2'', f_1) > (1-\epsilon)\Pi_2(f_2', f_1') + \epsilon\Pi_2(f_2''', f_1').$$

If the left and the right hand sides of the above inequality are equal, and if $f_2 = f_2''$ and $f_2' \neq f_2'''$, then player 2 prefers the former to the latter. Otherwise, player 2 is indifferent between the two outcomes. That is, if the supergame payoffs are the same, player 2 prefers to use the same supergame strategy chosen in the first stage rather than change his supergame strategy to another strategy. This implies that if the revision node for player 2 is reached and if the supergame strategy chosen by player 2 in the first stage maximizes his supergame payoff, then he strictly prefers his original supergame strategy chosen in the first stage to any other alternative.

3. ILLUSTRATION IN INFINITELY REPEATED PRISONERS' DILEMMA

In this section, we illustrate the logic of the main result in the repeated prisoners' dilemma. Figure 1 describes the one-shot prisoners' dilemma, in which (D, D) is the only Nash equilibrium. In the infinitely

repeated game with the limit of means criterion, every feasible pair of individually rational payoffs is supported by a subgame perfect equilibrium (the Folk Theorem).

Figure 1

In contrast, the set of pairs of equilibrium payoffs in this model consists of strictly individually rational and Pareto efficient allocations.^{5/} In other words, a pair of payoffs (x_1^*, x_2^*) is supported by a subgame perfect equilibrium if and only if either

$$(x_1^*, x_2^*) = \lambda(2,2) + (1-\lambda)(3,-1), \quad \lambda \in [\frac{1}{3}, 1],$$

or

$$(x_1^*, x_2^*) = \lambda(2,2) + (1-\lambda)(-1,3), \quad \lambda \in (\frac{1}{3}, 1],$$

holds. Particularly, (D, D) disappears as an equilibrium in this model. Figure 2 shows this.^{6/}

Figure 2

The logic goes as follows. Suppose that the both players use the same supergame strategy $f = \{f^t\}_{t=1}^{\infty}$ where

$$f^t(\cdot) = D \quad \text{for all } t.$$

Since there is a small probability that player 2 has a chance to revise his supergame strategy after being informed of player 1's choice, player 1 has an incentive to deviate (for example) to f_1 where

$$\begin{aligned} f_1^1 &= D, \\ f_1^t(s_1^1, \dots, s_2^{t-1}) &= \begin{cases} C & \text{if } s_2^{t-1} = C, \\ D & \text{if } s_2^{t-1} = D \text{ for } t > 1. \end{cases} \end{aligned} \quad \text{7/}$$

If player 1 deviates, and if the revision node for player 2 is reached, player 2 will choose (for example) f_2 where

$$f_2^1 = C,$$

$$f_2^t(\cdot) = c \text{ for } t > 1.$$

By this choice, both players get 2 instead of zero if the revision node for player 2 is reached. Knowing this, player 1 deviates to gain 2ε . Note that this new strategy profile is not an equilibrium, either, since player 2 deviates to f_2 not only in the third stage but also in the first stage. In the similar way, every Pareto dominated pair of payoffs does not sustain as a subgame perfect equilibrium in the leakage game.

4. EQUILIBRIA

In this section, we define subgame perfect equilibria of the leakage game and find the set of equilibrium outcomes. Given f_j , let $\bar{F}_i(f_j)$ be the class of the best-response-supergame-strategies, i.e., the class of all the supergame strategies f_i such that

$$\Pi_i(f_i, f_j) \geq \Pi_i(f'_i, f_j) \text{ for all } f'_i \in \mathcal{F}_i.$$

The mapping $\bar{F}_i: \mathcal{F}_j \rightarrow \mathcal{F}_i$ is the best response correspondence of player i to player j 's supergame strategies. Taking into account the cost of revision, we define the following. Given a pair of supergame strategies chosen in the first stage, let $\bar{F}_2(f_1, f_2)$ be:

$$\bar{F}_2(f_1, f_2) = \begin{cases} \{f_2\} & \text{if } f_2 \in \bar{F}_2(f_1) \\ \bar{F}_2(f_1) & \text{if } f_2 \notin \bar{F}_2(f_1). \end{cases}$$

$\bar{F}_2(f_1, f_2)$ is the best reaction correspondence in the third stage to f_1 given his choice f_2 in the first stage. If a pair of supergame strategies (f_1, f_2) is taken in the first stage, and if f_2 is the best response to f_1 , then player 2 strictly prefers to keep his supergame strategy rather than change it. Using this notion, we define the criterion for equilibria to the

leakage game.

Definition: Given ϵ , a strategy configuration $(f_1^*, f_2^*; F^*)$ is a subgame perfect equilibrium of the leakage game if

$$f_2^* \in \bar{F}_2(f_1^*), \quad (1)$$

$$F^*(f_1, f_2) \in \bar{F}_2(f_1, f_2) \quad \text{for all } f_1 \in \mathcal{F}_1 \text{ and } f_2 \in \mathcal{F}_2, \quad (2)$$

and for all $f_1' \in \mathcal{F}_1$

$$(1-\epsilon)\Pi_1(f_1^*, f_2^*) + \epsilon\Pi_1(f_1^*, F^*(f_1^*, f_2^*)) \geq (1-\epsilon)\Pi_1(f_1', f_2^*) + \epsilon\Pi_1(f_1', F^*(f_1', f_2^*)). \quad (3)$$

(2) implies that a best-reaction-supergame-strategy is chosen in every subgame of 'the leakage game. A function F^* satisfying (2) is called a best reaction function. The criterion for equilibria used in this paper is not subgame perfectness in the supergame; rather, it is perfectness in the leakage game, i.e., supergame strategy chosen in the third stage of the leakage game is the best reaction to the other supergame strategy given what was chosen in the first stage.

We are now in a position to state the following lemma which will be used later in proofs of theorems.

LEMMA: For any feasible pair of strictly individually rational payoffs (x_1^*, x_2^*) there exists a pair of supergame strategies (f_1, f_2) and a positive number d such that for $i=1,2$,

$$\begin{aligned} x_i^* &= \Pi_i(f_i, f_j) \\ f_i &\in \bar{F}_i(f_j), \end{aligned} \quad (4)$$

and

$$\Pi_i(f_i, f_j) \leq \Pi_i(f_i, f_j') \quad \text{for all } f_j' \text{ with } \Pi_j(f_j, f_i) - \Pi_j(f_j', f_i) < d. \quad (5)$$

Proof: Let $\{(s_1^t, s_2^t)\}_{t=1}^\infty$ be a sequence of pairs of actions such that

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \pi_i(s_1^t, s_2^t) = x_i^* \quad \text{for } i=1,2.$$

Next, we construct supergame strategies f_i ($i=1,2$) as follows:

$$f_i^1 = s_i^1,$$

$$f_i^t(r^1, \dots, r^{t-1}) = \begin{cases} s_i^t & \text{if } (r^1, \dots, r^{t-1}) = (s^1, \dots, s^{t-1}), \\ s_i & \text{otherwise.} \end{cases}$$

If player i tries to change the sequence of payoffs $((s_1^t, s_2^t))_{t=1}^\infty$, he can get at most π_i . Set $d = \min_{i=1,2} \{\frac{1}{2}(\Pi_i(f_i, f_j) - \pi_i)\}$. Then each supergame strategy is a best response to the other strategy, and each player cannot lower the supergame payoff of the other player without decreasing his own payoff by more than d . Q.E.D.

Since s_i is taken forever by player 1 if player j does not follow the strategies (s_j^t) , player j cannot change the outcome without lowering his payoff to π_j . Note that this lemma also shows that if a leakage game consists of only the first stage, i.e., if there is no possibility of revision of supergame strategy, then the Folk Theorem holds.

Before we present theorems, the definition of Pareto efficiency is given.

Definition: A feasible pair of payoffs (x_1^*, x_2^*) is (weakly) Pareto efficient if there exists no feasible pair of payoffs (x_1, x_2) such that

$$x_i > x_i^* \text{ for } i=1,2.$$

We say that a pair of payoffs is Pareto dominated if it is not Pareto efficient.

We are now in a position to characterize the set of pairs of payoffs that can be supported by subgame perfect equilibria. Theorem 1 is the main result of this paper stating that no Pareto dominated allocation can be supported by a subgame perfect equilibrium.

THEOREM 1: If $(\Pi_1(f_1, f_2), \Pi_2(f_2, f_1))$ is Pareto dominated, then for any best

reaction function F , $(f_1, f_2; F)$ is not a subgame perfect equilibrium of the leakage game.

Roughly speaking, the proof of the theorem goes as follows. Let (f'_1, f'_2) be a pair of supergame strategies satisfying $\Pi_i(f'_i, f'_j) > \Pi_i(f_i, f_j)$ for all $i=1,2$. Then both players prefer to change their supergame strategies from (f_1, f_2) to (f'_1, f'_2) if possible. Player 1 knows that if he constructs f'_1 , then with probability ϵ player 2 will revise his supergame strategy, and this revision will necessarily increase player 1's payoff as well as player 2's payoff. Here, however, the replacement of f_1 by f'_1 may harm player 1 because with probability of at least $1-\epsilon$ player 1 faces f_2 , and (f'_1, f_2) may decrease player 1's payoff. Therefore, player 1 constructs a supergame strategy f^*_1 by connecting f_1 and f'_1 ; his strategy f^*_1 behaves like f_1 unless it receives a certain signal from player 2 at stage 1 of the supergame, in which case f^*_1 immitates f'_1 . By using f^*_1 , (f^*_1, f_2) attains the same payoff as (f_1, f_2) does. If the revision node for player 2 is reached, player 2 now has an incentive to revise his supergame strategy to a strategy f^*_2 which gives f^*_1 the appropriate "signal" at stage 1 and immitates f'_2 after that. Therefore, player 1 has an incentive to construct f^*_1 . Note that the resulting strategy profile is not an equilibrium, either, since player 2 has an incentive to change his supergame strategy to f^*_2 not only in the third stage but also in the first stage. The following is the formal proof of the theorem.

Proof: Assume the contrary, i.e., that there exists an F such that $(f_1, f_2; F)$ is a subgame perfect equilibrium. First, note that $f_2 \in \bar{F}_2(f_1)$. For if not, player 2 has an incentive to deviate in the first stage. From the lemma, there exists a pair of supergame strategies (f'_1, f'_2) which

satisfies the following:

$$f'_i \in \bar{F}_i(f'_j) \quad (6)$$

$$\Pi_i(f'_i, f'_j) > \Pi_i(f_i, f_j) \quad (7)$$

and

$$\Pi_i(f'_i, f''_j) \geq \Pi_i(f'_i, f'_j) \text{ for all } f''_j \in \bar{F}_j(f'_i) \quad (8)$$

for $i=1,2$. We now construct a supergame strategy f_1^* by connecting f_1 and f'_1 in the following manner:

$$\begin{aligned} f_1^{*1} &= f_1^1, \\ f_1^{*t}(s^1, \dots, s^{t-1}) &= \begin{cases} f_1^t(s^1, \dots, s^{t-1}) & \text{if } s_2^1 = f_2^1, \\ f_1^{t-1}(s^2, \dots, s^{t-1}) & \text{if } s_2^1 \neq f_2^1 \text{ for } t > 1. \end{cases} \end{aligned}$$

If the node in which player 2 can revise his supergame strategy is reached, player 2 has an incentive to change his supergame strategy to f_2^* , which is constructed as follows:

$$\begin{aligned} f_2^{*1} &\neq f_2^1, \\ f_2^{*t}(s^1, \dots, s^{t-1}) &= f_2^{t-1}(s^2, \dots, s^{t-1}) \text{ for } t > 1. \end{aligned}$$

Clearly, $f_2^* \in \bar{F}_2(f_1^*)$, and (f_1^*, f_2^*) attains the same pair of payoffs as (f'_1, f'_2) does. Let F be any best reaction function. Then

$$\Pi_1(f_1^*, F(f_1^*, f_2)) \geq \Pi_1(f'_1, f'_2)$$

always holds by virtue of (8). Furthermore, $F(f_1, f_2) = f_2$ holds since $\bar{F}_2(f_1, f_2) = \{f_2\}$ by $f_2 \in \bar{F}_2(f_1)$. Since we have

$$\Pi_1(f_1^*, f_2) = \Pi_1(f_1, f_2)$$

and

$$\Pi_1(f_1^*, f_2^*) = \Pi_1(f'_1, f'_2) > \Pi_1(f_1, f_2),$$

the following inequality holds:

$$(1-\epsilon)\Pi_1(f_1^*, f_2) + \epsilon\Pi_1(f_1^*, F(f_1^*, f_2)) > (1-\epsilon)\Pi_1(f_1, f_2) + \epsilon\Pi_1(f_1, F(f_1, f_2))$$

for any best reaction function F . Therefore, player 1 has an incentive to

deviate from f_1 to f_1^* for any best reaction function F . Hence, this is not an equilibrium. Q.E.D.

Since player 1 usually faces the supergame strategy f_2 , he constructs f_1^* to cope with that supergame strategy, too. Note that the above logic works because $\bar{F}_2(f_1, f_2)$ is a singleton. In other words, f_2 is not only a best-reaction but also the best-reaction-supergame-strategy to f_1^* given f_2 . This shows that the revision cost as well as the possibility of revision is essential to the result.

Indeed, if there is no cost of revision, then the following strategy profile $(f_1^*, f_2^*; F^*)$ of the information leakage game is an equilibrium in the example of Section 3:

$$f_1^{*t}(\cdot) = \begin{cases} C & \text{if } s_2^1, \dots, s_2^{t-1} = C, \text{ and } t = 3m \text{ where } m = 1, 2, \dots, \\ D & \text{otherwise,} \end{cases}$$

$$f_2^{*t}(\cdot) = D \quad \text{for all } t,$$

and F^* is a best response function, i.e., $F^*(f_1, f_2) \in \bar{F}_2(f_1)$ for all $f_1 \in F_1$ and all $f_2 \in F_2$, with

$$F^*(f_1^*, f_2^*)^t(\cdot) = C \quad \text{for all } t.$$

Here, player 1 takes the strategy that cooperates every three periods as long as player 2 has cooperated. First, player 2 has no incentive to deviate since F^* is a best response function, and f_2^* is a best response to f_1^* . Since f_1^* is a best response to f_2^* , player 1 has an incentive to deviate, say, to f_1' only if $\Pi_1(f_1', F^*(f_1', f_2^*))$ is greater than $\Pi(f_1^*, F^*(f_1^*, f_2^*)) = 2\frac{2}{3}$. If this is the case, then player 2's payoff should be below the individually rational payoff, which contradicts that F^* is a best response function. Thus player 1 has no incentive to deviate.

Next, we derive the following theorem which states that every strictly

individually rational and Pareto efficient allocation is supported by a subgame perfect equilibrium.

THEOREM 2: Suppose a feasible pair of payoffs (x_1^*, x_2^*) is Pareto efficient and strictly individually rational. Then, there exists a triple $(f_1^*, f_2^*; F^*)$ and $\bar{\epsilon} > 0$ such that for all $0 < \epsilon \leq \bar{\epsilon}$,

$$\Pi_i(f_1^*, f_2^*) = x_i^* \quad i=1,2,$$

and $(f_1^*, f_2^*; F^*)$ is a subgame perfect equilibrium with respect to ϵ .

The logic of the proof of the theorem goes as follows. First, from the lemma, we can construct a pair of supergame strategies (f_1^*, f_2^*) satisfying (4) and (5) in which (f_1, f_2) is replaced by (f_1^*, f_2^*) . Player 2 cannot benefit from unilateral deviation. As for player 1, he cannot gain anything either unless player 2 responds to his deviation in the revision node. Consider the extreme case in which the revision node is reached almost certainly, i.e., $\epsilon \approx 1$. In this situation player 1 can "blackmail" player 2 by choosing a strategy which makes player 2 play the strategy that gives player 1 the payoff higher than x_1^* . If the possibility of reaching revision node is small, however, player 1 should also take into account the possibility that his blackmail cannot affect player 2's supergame strategy, and his strategy has to face f_2^* . In that case, the probability of which is $1-\epsilon$, player 1's payoff decreases by at least d ; recall that (f_1^*, f_2^*) is constructed so that each player cannot decrease his opponent's payoff without decreasing his own payoff by at least d . For sufficiently small ϵ , this loss cannot be compensated by the gain obtained through the revision of supergame strategy by player 2. Therefore, player 1 has no incentive to deviate. The following is the formal proof of the theorem.

Proof: From the lemma, there exists a pair of supergame strategies (f_1^*, f_2^*)

such that for $i=1,2$,

$$\Pi_i(f_i^*, f_j^*) = x_i^*, \quad (9)$$

$$f_i^* \in \bar{F}_i(f_j^*), \quad (10)$$

$$\Pi_i(f_i^*, f_j^*) \leq \Pi_i(f_i^*, f_j') \quad \text{for all } f_j' \text{ with } \Pi_j(f_i^*, f_j^*) - \Pi_j(f_i^*, f_j') < d$$

for some $d > 0$, (11)

(11) implies that neither of the two players can lower the other's payoff without decreasing his own payoff by at least d . Let $F^*(f_1, f_2)$ be in $\bar{F}_2(f_1, f_2)$ for all $f_1 \in F_1$ and all $f_2 \in F_2$. We will prove that $(f_1^*, f_2^*; F^*)$ is a subgame perfect equilibrium of the information leakage game for all $0 < \epsilon \leq \bar{\epsilon}$ for some $\bar{\epsilon} > 0$. Clearly, F^* is a best reaction function, and player 2 has an incentive to deviate neither in the first stage nor in the third stage since (10) holds. Therefore, the triple $(f_1^*, f_2^*; F^*)$ sustains itself as a subgame perfect equilibrium if and only if for all $f_1'' \in F_1$,

$$(1-\epsilon)\Pi_1(f_1^*, f_2^*) + \epsilon\Pi_1(f_1^*, F^*(f_1^*, f_2^*))$$

$$\geq (1-\epsilon)\Pi_1(f_1'', f_2^*) + \epsilon\Pi_1(f_1'', F^*(f_1'', f_2^*)) \quad (12)$$

holds. Let then f_1'' be an arbitrary strategy in F_1 . To show that (12) holds, we consider two cases.

i) If $\Pi_1(f_1'', f_2^*) = \Pi_1(f_1^*, f_2^*)$, i.e., $f_1'' \in \bar{F}_1(f_2^*)$, then (11) implies that

$$\Pi_2(f_2^*, f_1'') \geq \Pi_2(f_2^*, f_1^*).$$

We consider the following two subcases.

i-a) If $f_2^* \notin \bar{F}_2(f_1'')$, then $F^*(f_1'', f_2^*) \in \bar{F}_2(f_1'')$ implies that

$$\Pi_2(F^*(f_1'', f_2^*), f_1'') > \Pi_2(f_2^*, f_1'') \geq \Pi_2(f_2^*, f_1^*).$$

From the Pareto optimality of $(\Pi_1(f_1^*, f_2^*), \Pi_2(f_2^*, f_1^*))$, we have

$$\Pi_1(f_1'', F^*(f_1'', f_2^*)) \leq \Pi_1(f_1^*, f_2^*).$$

Therefore, (12) holds for all ϵ .

i-b) If $f_2^* \in \bar{F}_2(f_1'')$, then player 2 does not have an incentive to revise his

supergame strategy even when he has a chance to do so in the third stage, i.e., $F^*(f_1^'', f_2^*) = f_2^*$, and (12) trivially holds for all ϵ .

ii) Next, suppose $\Pi_1(f_1^'', f_2^*) < \Pi_1(f_1^*, f_2^*)$. Then to violate (12),

$$(1-\epsilon)(\Pi_1(f_1^*, f_2^*) - \Pi_1(f_1^'', f_2^*)) < \epsilon(\Pi_1(f_1^'', F^*(f_1^'', f_2^*)) - \Pi_1(f_1^*, f_2^*)) \quad (13)$$

should hold. If $f_2^* \in \bar{F}_2(f_1^'')$, then $F^*(f_1^'', f_2^*) = f_2^*$, and (13) fails to hold by $f_1^* \in \bar{F}_1(f_2^*)$. Hence,

$$f_2^* \notin \bar{F}_2(f_1^'') \quad (14)$$

holds. (14), the positiveness of the right hand side of (13), and the Pareto optimality imply that

$$\Pi_2(f_2^*, f_1^'') < \Pi_2(F^*(f_1^'', f_2^*), f_1^'') \leq \Pi_2(f_2^*, f_1^*).$$

From (11), it must be the case that

$$\Pi_1(f_1^*, f_2^*) - \Pi_1(f_1^'', f_2^*) > d.$$

Therefore, for sufficiently small ϵ (13) does not hold since the set of feasible payoffs is bounded. Hence (12) holds for sufficiently small ϵ .

Since $f_1^''$ is arbitrary, player 1 has no incentive to deviate. Therefore, $(f_1^*, f_2^*; F^*)$ is a subgame perfect equilibrium. Q.E.D.

5. EXTENSIONS OF THE RESULT

The result of the previous section can be extended to other situations. Here, we examine two of them. The first is that both players have a chance to revise their supergame strategies. The second situation is that revision cost is small but positive so that preference relation is not lexicographic. In the following only the counterparts of Theorem 1 are presented.

(1) Two-sided espionage

In this subsection the information leakage game is modified to describe the situation in which both players have a chance to be informed of the other players' supergame strategies. The structure of the leakage game is the same as before except that player 1 as well as player 2 has a chance to revise his strategy after observing the other player's strategy with probability ϵ . It is assumed that both players' information gathering activities are independent of each other without knowing whether their own strategies are revealed to their opponents. Therefore, they simultaneously have chances to revise their strategies with probability ϵ^2 . A strategy profile is written as $(f_1, f_2; F_1, F_2)$ where $f_i \in \mathcal{F}_i$ and $F_i: \mathcal{F}_i \times \mathcal{F}_j \rightarrow \mathcal{F}_i$ for $i=1,2$. The expected supergame payoff for player i ($i=1,2$) is now given by

$$\begin{aligned} \bar{\Pi}_i(f_i, f_j; F_i, F_j) = & (1-\epsilon)^2 \Pi_i(f_i, f_j) + \epsilon(1-\epsilon) \Pi_i(F_i(f_i, f_j), f_j) \\ & + \epsilon(1-\epsilon) \Pi_i(f_i, F_j(f_j, f_i)) + \epsilon^2 \Pi_i(F_i(f_i, f_j), F_j(f_j, f_i)). \end{aligned}$$

A similar modification is made on preference relations. We assume that player 1 has lexicographic preference as player 2 does, that is to say, the expected supergame payoff for player 1 matters first, and whether player 1 changes his supergame strategy in the revision node matters second.

We apply subgame perfect equilibria to this modified leakage game. A strategy profile $(f_1^*, f_2^*; F_1^*, F_2^*)$ is said to be a subgame perfect equilibrium if

$$\begin{aligned} \bar{\Pi}_i(f_i^*, f_j^*; F_i^*, F_j^*) \geq \bar{\Pi}_i(f_i, f_j^*; F_i, F_j^*) \quad & \text{for all } f_i \in \mathcal{F}_i \text{ and all } F_i, \text{ and} \\ F_i^*(f_i, f_j) \in \operatorname{argmax}_{f_i' \in \mathcal{F}_i} & (1-\epsilon) \Pi_i(f_i', f_j) + \epsilon \Pi_i(f_i', F_j^*(f_j, f_i)) \\ & \text{for all } f_i \in \mathcal{F}_i \text{ and all } f_j \in \mathcal{F}_j, \end{aligned}$$

with

$$F_i^*(f_i, f_j) = f_i \quad \text{if } f_i \in \operatorname{argmax}_{f_i' \in \mathcal{F}_i} (1-\epsilon) \Pi_i(f_i', f_j) + \epsilon \Pi_i(f_i', F_j^*(f_j, f_i)),$$

for $i=1,2$.

If both players have large chances of observing other player's supergame strategies, for example with probability one, the result will be the same as in the case of no information leakage. To see this in the example of Section 3, suppose that both players take the supergame strategy that always defects in the first stage, and suppose that both players' reaction functions are such that each player chooses the same strategy when he takes the supergame strategy that always defects in the first stage. Then this strategy profile becomes a subgame perfect equilibrium for properly chosen F_1 and F_2 . This happens because each player's observation has nothing to do with the actual supergame strategy taken by his opponent since it is revised with probability one. Thus, we have the following weaker statement the proof of which is given in the appendix.

THEOREM 3: If for any $\bar{\epsilon} > 0$ there exists $\epsilon \in (0, \bar{\epsilon}]$ such that $(f_1, f_2; F_1, F_2)$ is a subgame perfect equilibrium of the leakage game with respect to ϵ , then $(\Pi_1(f_1, f_2), \Pi_2(f_2, f_1))$ is Pareto efficient.

Proof: See Appendix.

(2) Positive revision cost

This subsection examines the case when the cost of revision is large so that the loss in supergame payoff may be compensated by not revising his strategy. To make matters simple we turn back to the case of one-sided information leakage. We assume that the cost of revision incurred by player 2 is $k (> 0)$. Since we are no longer bothered by lexicographical preference, we define the total payoff for player 2 as follows:

$$\bar{\Pi}_2(f_2, f_1; F) = (1-\epsilon)\Pi_2(f_2, f_1) + \epsilon(\Pi_2(F(f_1, f_2), f_1) - k\delta)$$

where $\delta = \begin{cases} 0 & \text{if } F(f_1, f_2) = f_2, \\ 1 & \text{otherwise.} \end{cases}$

Player 1's total payoff is defined in a similar way except that $k=0$ for player 1. Note that the revision cost is incurred by player 2 only when he actually changes his supergame strategy. Then a strategy profile $(f_1^*, f_2^*; F^*)$ is a subgame perfect equilibrium of the leakage game if

$$\begin{aligned} \bar{\Pi}_1(f_1^*, f_2^*; F^*) &\geq \bar{\Pi}_1(f_1, f_2^*; F^*) && \text{for all } f_1 \in \mathcal{F}_1, \\ \bar{\Pi}_2(f_2^*, f_1^*; F^*) &\geq \bar{\Pi}_2(f_2, f_1^*; F) && \text{for all } f_2 \in \mathcal{F}_2 \text{ and all } F, \text{ and} \\ \bar{\Pi}_2(f_2', f_1'; F^*) &\geq \bar{\Pi}_2(f_2', f_1'; F) && \text{for all } F, \text{ for any } f_1' \in \mathcal{F}_1 \text{ and } f_2' \in \mathcal{F}_2. \end{aligned}$$

Suppose that both players' supergame strategies chosen in the first stage are best responses to each other. Then player 1 gains by deviation only if player 2 responds to it in the revision node, that is to say, only if player 2's supergame payoff is increased by more than k ; still player 1 should get more than his current supergame payoff. Thus, the statement corresponding to Theorem 1 is modified in the following way.

THEOREM 4: If $(\Pi_1(f_1, f_2), \Pi_2(f_2, f_1) + k)$ is Pareto dominated, Then $(f_1, f_2; F)$ cannot be a subgame perfect equilibrium of the leakage game for any $F: \mathcal{F}_1 \times \mathcal{F}_2 \rightarrow \mathcal{F}_2$.

We omit the proof of the theorem since its logic is the same as that of Theorem 1.

6. CONCLUSION

In a supergame without information leakage, Pareto dominated outcomes are supported by equilibria because each player believes that his deviation cannot affect the strategy of his opponent. This sense of resignation is

conducive to the appearance of a Pareto dominated outcome as an equilibrium. If a player knows that his deviation leads to revision of strategy of his opponent with at least small probability, the player incorporates the following statement into his strategy: if his opponent behaves cooperatively, then he will also cooperate. If this is indeed the case, then his opponent follows cooperative behavior. In this manner, Pareto dominated outcomes cannot be stable, and cooperation emerges.

There are two remarks on this model. First, in this model the revision of strategy is assumed to occur before the supergame starts. This revision may occur at any stage of the supergame as far as the total probability of revision is small enough.

Second, I constructed an information leakage game and applied Nash (subgame perfect) equilibria to this game. There is an underlying assumption that both players determine their entire strategy before the game begins. This assumption does not hold in the real world; rather, people make trial and error to seek the best (or better) strategy for them.^{8/} The question is what will occur under such circumstances. I would like to conclude the paper by raising this question not only as the problem of the model of this paper, but also as the problem of the current approaches to the repeated games.

FOOTNOTES

- 1) See the survey by Aumann [4].
- 2) As for finitely repeated games, a serious problem arises when one tries to relate it to infinitely repeated games. Selten [13] pointed out that no matter how long the horizon becomes, the equilibrium outcome is simply the repetition of the equilibrium outcome of a one-shot game unless the component game has a special structure.
- 3) In a pioneering paper Friedman [6] derived Pareto efficient outcomes as noncooperative equilibria, but he did so simply by assuming it.
- 4) To be honest, I have not seen the proof of the result obtained in Aumann and Sorin [5].
- 5) As will be defined in Section 4, the criterion for equilibria used in this model is subgame perfectness in the sense that every player has an incentive to change his supergame strategy neither in the first stage nor in the third stage.
- 6) This asymmetry at the end point is caused by the fact that only player 2 has a chance to revise the supergame strategy. Note that if the information on player 1's supergame strategy leaks to player 2 with probability 1, i.e., $\epsilon=1$, then the set of equilibrium pair of payoffs consists of $(2\frac{2}{3}, 0)$ only.
- 7) As you might see, this is not the best deviation for player 1.
- 8) Kaneko [8] proposes an alternative framework and solution concept called conventionally stable sets. There, players do not know the structure of the game and try to attain higher stationary payoff by trial and error.

REFERENCES

- [1] Abreu, D., "Extremal Equilibria of Oligopolistic Supergames," Journal of Economic Theory, 39, 1986, pp.191-225.
- [2] Abreu, D., and A. Rubinstein, "The Structure of Nash Equilibrium in Repeated Games with Finite Automata," mimeo, 1986.
- [3] Aumann, R., "Acceptable Points in General Cooperative n-Person Games," in Contributions to the Theory of Games IV eds. by A.W. Tucker and R.D. Luce, 1959, Princeton University Press.
- [4] Aumann, R., "Survey of Repeated Games," in Essays in Game Theory and Mathematical Economics in Honor of Oskar Morgenstern, Bibliographisches Institut Mannheim,Wien,Zurich.
- [5] Aumann, R., and S. Sorin, "Cooperation and Bounded Rationality," mimeo, 1986.
- [6] Friedman, J.W., "A Non-cooperative Equilibrium for Supergames," Review of Economic Studies, 38, 1971, pp.1-12.
- [7] Kalai, E., "Preplay Negotiations and the Prisoner's Dilemma," Mathematical Social Sciences, 1, 1981, pp.375-379.
- [8] Kaneko, M., "The Conventionally Stable Sets in Noncooperative Games with Limited Observations I: Definitions and Introductory Arguments," Mathematical Social Sciences, 13, 1987, pp.93-128.
- [9] Kreps, D., P. Milgrom, J. Roberts, and R. Wilson, "Rational Cooperation in the Repeated Prisoners' Dilemma," Journal of Economic Theory, 27, 1982, pp.245-252.
- [10] Pearce, D.G., "Renegotiation-Proof Equilibria: Collective Rationality

and Intertemporal Cooperation," mimeo, 1987.

- [11] Rubinstein, A., "Equilibrium in Supergames with the Overtaking Criterion," Journal of Economic Theory, 21, 1979, pp.1-9.
- [12] Rubinstein, A., "Finite Automata Play the Repeated Prisoners' Dilemma," Journal of Economic Theory, 39, 1986, pp.83-96.
- [13] Selten, R., "The Chain Store Paradox," Theory and Decision, 9, 1978, pp.127-159.

APPENDIX

In this appendix, we present the proof of Theorem 3. To do that, we need the following two lemmas.

LEMMA A1: Suppose that there exists $(f'_1, f'_2) \in \mathcal{F}_1 \times \mathcal{F}_2$ such that

$$\Pi_i(f'_i, f'_j) > \Pi_i(f_i, F_j(f_j, f_i)) \quad (\text{a-1})$$

for $i=1,2$. Then $(f_1, f_2; F_1, F_2)$ is not a subgame perfect equilibrium for any $\epsilon \in (0, \bar{\epsilon}]$ for some $\bar{\epsilon} > 0$.

Proof: Assume the contrary, i.e., that there exist F_1 and F_2 such that $(f_1, f_2; F_1, F_2)$ is a subgame perfect equilibrium for $0 < \epsilon \leq \bar{\epsilon}$. From the lemma in the main text, there exists a pair of supergame strategies $(g_1, g_2) \in \mathcal{F}_1 \times \mathcal{F}_2$ which satisfies the following:

$$g_i \in \bar{F}_i(g_j), \quad (\text{a-2})$$

$$\Pi_i(g_i, g_j) = \Pi_i(f'_i, f'_j), \quad (\text{a-3})$$

and

$$\begin{aligned} \Pi_i(g_i, g'_j) &\geq \Pi_i(g_i, g_j) \text{ for all } g'_j \in \mathcal{F}_j \text{ satisfying} \\ &\Pi_j(g_j, g_i) - \Pi_j(g'_j, g_i) < d \text{ for some } d > 0 \end{aligned} \quad (\text{a-4})$$

for $i=1,2$. Construct a supergame strategy f_1^* by connecting g_1 and f_1 in the following manner:

$$\begin{aligned} f_1^{*1} &= f_1^1, \\ f_1^{*t}(s^1, \dots, s^{t-1}) &= \begin{cases} f_1^t(s^1, \dots, s^{t-1}) & \text{if } s_2^1 = f_2^1, \\ g_1^{t-1}(s^2, \dots, s^{t-1}) & \text{if } s_2^1 \neq f_2^1 \text{ for } t > 1. \end{cases} \end{aligned}$$

From (a-4), for sufficiently small $\epsilon > 0$

$$\Pi_1(f_1^*, F_2(f_2, f_1^*)) \geq \Pi_1(g_1, g_2)$$

holds. Moreover, we have

$$\Pi_1(f_1^*, f_2) = \Pi_1(f_1, f_2),$$

$$\Pi_1(F_1(f_1^*, f_2), f_2) - \Pi_1(F_1(f_1, f_2), f_2).$$

Then

$$\begin{aligned} \bar{\Pi}_1(f_1^*, f_2; F_1, F_2) &= (1-\varepsilon)^2 \Pi_1(f_1^*, f_2) + (1-\varepsilon)\varepsilon \Pi_1(f_1^*, F_2(f_2, f_1^*)) \\ &\quad + (1-\varepsilon)\varepsilon \Pi_1(F_1(f_1^*, f_2), f_2) + \varepsilon^2 \Pi_1(\cdot) \\ &> (1-\varepsilon)^2 \Pi_1(f_1, f_2) + (1-\varepsilon)\varepsilon \Pi_1(f_1, F_2(f_2, f_1)) \\ &\quad + (1-\varepsilon)\varepsilon \Pi_1(F_1(f_1, f_2), f_2) + \varepsilon^2 \Pi_1(\cdot) \\ &= \bar{\Pi}_1(f_1, f_2; F_1, F_2) \end{aligned}$$

holds for sufficiently small ε by virtue of (a-1) and (a-3). Thus, player 1 has an incentive to deviate from f_1 to f_1^* in the first stage, and then $(f_1, f_2; F_1, F_2)$ is not a subgame perfect equilibrium for sufficiently small $\varepsilon > 0$. Q.E.D.

LEMMA A2: Suppose that for any $\bar{\varepsilon} > 0$ there exists $\varepsilon \in (0, \bar{\varepsilon}]$ such that $(f_1, f_2; F_1, F_2)$ is a subgame perfect equilibrium with respect to ε . Then $F_1(f_1, f_2) \neq f_1$ implies $F_2(f_2, f_1) \neq f_2$ and

$$\begin{aligned} \Pi_i(F_i(f_i, f_j), f_j) &= \Pi_i(f_i, f_j) \\ \Pi_i(F_i(f_i, f_j), F_j(f_j, f_i)) &> \Pi_i(f_i, F_j(f_j, f_i)) \end{aligned}$$

for $i=1, 2$.

Proof: Let $g_i = F_i(f_i, f_j)$ for $i=1, 2$. The subgame perfectness of $(f_1, f_2; F_1, F_2)$ and $g_1 \neq f_1$ implies that

$$(1-\varepsilon)\Pi_1(g_1, f_2) + \varepsilon\Pi_1(g_1, g_2) > (1-\varepsilon)\Pi_1(f_1, f_2) + \varepsilon\Pi_1(f_1, g_2).$$

This inequality must hold for sufficiently small ε . Thus, we have

$$\Pi_1(g_1, f_2) \geq \Pi_1(f_1, f_2) \tag{a-5}$$

where the equality holds only if

$$\Pi_1(g_1, g_2) > \Pi_1(f_1, g_2).$$

If (a-5) holds with strict inequality, then player 1 has an incentive to

deviate from f_1 to g_1 in the first stage for sufficiently small ϵ . Therefore it must be the case that

$$\Pi_1(g_1, f_2) = \Pi_1(f_1, f_2), \text{ and}$$

$$\Pi_1(g_1, g_2) > \Pi_1(f_1, g_2),$$

which implies that $g_2 \neq f_2$. Hence, in a similar way, we obtain

$$\Pi_2(g_2, f_1) = \Pi_2(f_2, f_1),$$

$$\Pi_2(g_2, g_1) > \Pi_2(f_2, g_1).$$

Q.E.D.

We are now in a position to present the proof of Theorem 3, which is rather straightforward once we prove the above two lemmas.

Proof of Theorem 3: Suppose $(f_1, f_2; F_1, F_2)$ is a subgame perfect equilibrium, and $(\Pi_1(f_1, f_2), \Pi_2(f_2, f_1))$ is Pareto dominated. If $F_i(f_i, f_j) = f_i$ for $i=1, 2$, then there exists $(f'_1, f'_2) \in \mathcal{F}_1 \times \mathcal{F}_2$ such that

$$\Pi_i(f'_i, f'_j) > \Pi_i(f_i, F_j(f_j, f_i)) = \Pi_i(f_i, f_j)$$

for $i=1, 2$. Then by Lemma A1 $(f_1, f_2; F_1, F_2)$ is not a subgame perfect equilibrium. So suppose $F_i(f_i, f_j) \neq f_i$ for $i=1$ or 2 . Let $g_i = F_i(f_i, f_j)$ for $i=1, 2$. By Lemma A2,

$$\Pi_i(g_i, g_j) > \Pi_i(f_i, F_j(f_j, f_i))$$

holds for $i=1, 2$. Again by Lemma A1, $(f_1, f_2; F_1, F_2)$ is not a subgame perfect equilibrium.

Q.E.D.

Player 2

		C	D
Player 1	C	2 2	-1 3
	D	3 -1	0 0

Figure 1

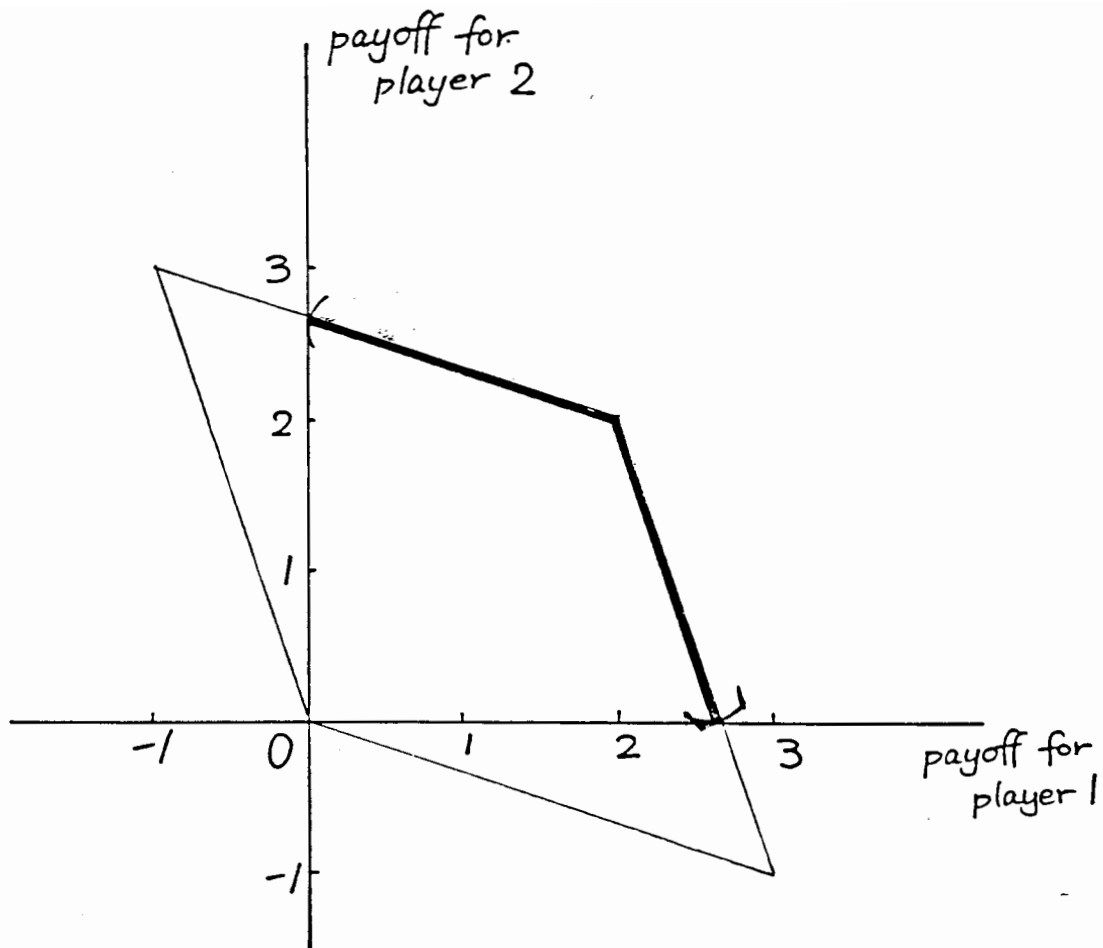


Figure 2