

Digitized by the Internet Archive
in 2011 with funding from
Boston Library Consortium Member Libraries

<http://www.archive.org/details/folktheoremstwod00smit>

DEWEY

MIT LIBRARIES



3 9080 00756965 7

HB31
.M415

no. 597

**working paper
department
of economics**

Folk Theorems: Two-Dimensionality is
(Almost) Enough

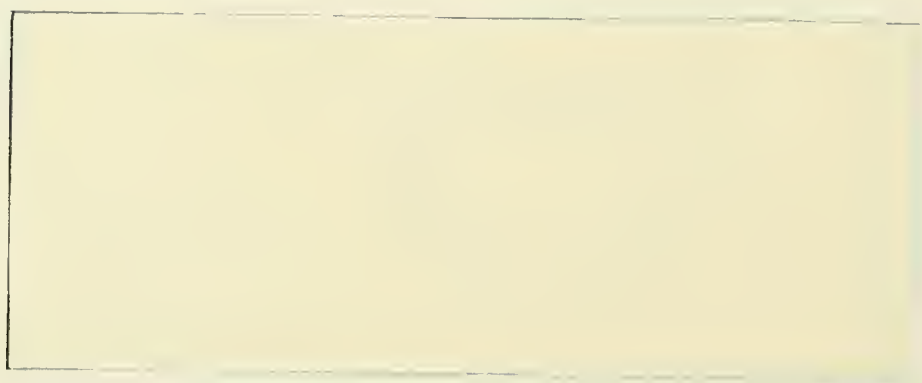
Lones Smith
Department of Economics
M.I.T.

No. 597

Jan. 1992

**massachusetts
institute of
technology**

**50 memorial drive
cambridge, mass. 02139**



Folk Theorems: Two-Dimensionality is
(Almost) Enough

Lones Smith
Department of Economics
M.I.T.

No. 597

Jan. 1992

M.I.T. LIBRARIES
MAR 09 1992
RECEIVED

Folk Theorems: Two-Dimensionality is (Almost) Enough*

Lones Smith
Department of Economics
M. I. T.

original version: October 8, 1990
current version: January 22, 1992

Abstract

We exhibit a parallel constructive proof technique for three foundational folk theorems *without* requiring a full-dimensional payoff space. Instead, we show that a less stringent and more natural projection condition suffices. Along the way, we substitute a strengthened folk theorem for finitely-repeated games.

*The current version reflects helpful comments from two referees. Financial assistance from the Social Sciences and Humanities Research Council of Canada, the Jacob K. Javits Fellowship Fund, and the Searle Foundation is gratefully acknowledged.

1. INTRODUCTION

The problem of multiplicity of equilibria in repeated games is by now well-known. Indeed, the popular moniker ‘folk theorem’ has been effectively hijacked by game theorists to mean that any individually rational payoff of a stage game can be attained on average in an equilibrium of the corresponding repeated game. Constructing n -player folk theorems for various conceivable economic contexts has developed into a veritable cottage industry. Most recent efforts can trace their lineage to the perfect folk theorem in Rubinstein (1979), which considered supergames without discounting. When players do discount the future, Rubinstein’s recourse to infinitely-bootstrapped punishment hierarchies is no longer an option. It was a technique first pioneered by Fudenberg and Maskin (1986) (henceforth simply F-M) for infinite horizon games that sparked the latest flurry in this arena. The requirement in F-M that the feasible payoff space have *full dimension* — i.e. a non-empty interior — found immediate application in the folk theorem (without discounting) for finitely repeated games in Benoit and Krishna (1985) (hereafter denoted B-K). Smith (1990) later exploited the full-dimensionality condition to establish a *uniform* folk theorem for overlapping generations games.

Exactly why a full-dimensional payoff space was deemed necessary is best explained by means of an example game G_1 found in B-K:

3, 3, 3	0, 0, 0	1, 1, 1	2, 2, 2
0, 1, 1	0, 0, 0	0, 1, 1	0, 0, 0
0, 0, 0	0, 0, 0	0, 1, 1	0, 1, 1

Here, player 1 chooses rows, 2 chooses columns, and 3 chooses matrices. There are two Nash equilibria: $(3, 3, 3)$ and $(2, 2, 2)$. Since each player’s minimax level is 0, a folk theorem would assert that any positive feasible payoff vector should be attainable (on average) in $G_1(T)$, the T -fold repetition of G_1 , for T sufficiently large.

Notice that the payoff space of G_1 is two-dimensional, and thus is not of full-dimension.¹ But the choice of G_1 was really rather clever! For consider the simple proof in B-K that in no subgame perfect equilibrium of any $G_1(T)$ does player 2 or 3 earn less than 1 on average. The result proceeds by (backward) induction. The case $T = 1$ is obvious, so assume it holds up to some $T \geq 1$. Then the prescribed outcome with $T + 1$ periods to go cannot be $(0, 0, 0)$, for if so, either player 2 or 3 could gain at least 1 by deviating and then, by hypothesis, average at least 1 in the following subgame. But if $(0, 0, 0)$ is not prescribed, then players 2 and 3 each receive at least 1, and the claim follows for $G_1(T + 1)$, as required.

Unfortunately, this neat result relies *crucially* on the fact that players 2 and 3 receive positively linearly related payoffs. It turns out that such a non-genericity,

¹The example in F-M only has a 1-dimensional payoff space, and hence might not establish that a 3-dimensional payoff space is needed.

which was conceded by B-K, is not without loss of generality. In the sequel, we find that when this possibility is expressly eliminated, the folk theorems for $n > 2$ players in F-M, B-K, and Smith (1990) obtain with any payoff space that is (at least) two-dimensional. This task actually proves rather unburdensome, as a unified constructive treatment suffices for all three cases. We also take the liberty of strengthening the folk theorem of B-K, by admitting the possibility that one player might not have distinct Nash payoffs in the stage game.

2. PUNISHMENT OUTCOMES IN 2-SPACE

The set-up is standard, and is taken from Smith (1990), which we summarize here. Throughout, G is an n -person normal form game, and the (compact and convex) strategy spaces $\{A_1, \dots, A_n\}$ consist of mixed strategies. Each player i 's payoff function $U_i : \prod_{j=1}^n A_j \rightarrow \mathcal{R}$ is continuous. Define $U = \prod_{j=1}^n U_j$. Elements $a^* \in A$ are referred to as *outcomes* of G . Liberal use is made of *correlated* strategies, allowing players to condition their actions on the outcome of a public randomizing device. This renders the *feasible* payoff space V the convex hull in \mathcal{R}^n of the set of pure strategy payoff vectors.

Let M^i be a minimax strategy against player i in the game G . Since M^i could be a mixed strategy, we assume that deviations from mixed strategies are observable.² We may also assume WLOG that $U_i(M^i) = 0$ for all $i \in \mathbf{N} \equiv \{1, 2, \dots, n\}$. The feasible and *strictly individually rational* payoff set is therefore the *positive* orthant V^* of V .

The example game G_1 was atypical in the following sense: Even though the whole payoff space was two-dimensional, its projection onto the plane of payoffs for players 2 and 3 was only an (upward-sloping) line. It was thus impossible to decouple the payoffs of 2 and 3. We claim that so long as we avoid such a non-genericity for n -player games, we can construct n *personalized one-shot (correlated) punishment outcomes* $p_1^*, p_2^*, \dots, p_n^*$. That is, (i) p_k^* is a strictly individually rational outcome, (ii) providing k a lower payoff than the equilibrium outcome a^* , and (iii) a strictly lower payoff than any other $p_i^*, i \neq k$. Observe that because we do not have a full-dimensional payoff space, other players $j \neq k$ may suffer from p_k^* ; however, every $j \neq k$ still prefers p_k^* to p_j^* .

For the sake of definiteness, we shall (locally) parameterize V^* assuming it is at most two-dimensional. The extension to higher dimensions should be obvious. Let the payoff of each player $i \in \mathbf{N}$ be $\pi_i(s, t) = \alpha_i + \beta_i s + \gamma_i t$. All subsequent analysis rests on the following crucial assumption which will supplant full-dimensionality:

- (P) The projection V_{ij} of V^* onto the coordinate space of any two players j and k is either two-dimensional, so that $\beta_j \gamma_k \neq \beta_k \gamma_j$ for all $j \neq k$, or a line with negative slope, i.e. $\beta_j = \theta_{jk} \beta_k$ and $\gamma_j = \theta_{jk} \gamma_k$ for some $\theta_{jk} < 0$.

For example, with $n = 3$ players, (P) implies that if the payoff space is perpendicular to the coordinate plane of any two players, then it intersects that plane in a

²Abreu-Dutta (1991) modify the argument of F-M to handle unobservable mixed strategies.

negatively sloped line.

Now suppose that we want to support the strictly individually rational outcome a^* , as parameterized by (s^*, t^*) . Under (P), we can exhibit coordinates $(s_1, t_1), \dots, (s_n, t_n)$ for $p_1^*, p_2^*, \dots, p_n^*$ satisfying conditions (i), (ii), and (iii).³ Let

$$\begin{aligned} s_i &= s^* - \beta_i \varepsilon / \sqrt{\beta_i^2 + \gamma_i^2} \\ t_i &= t^* - \gamma_i \varepsilon / \sqrt{\beta_i^2 + \gamma_i^2} \end{aligned}$$

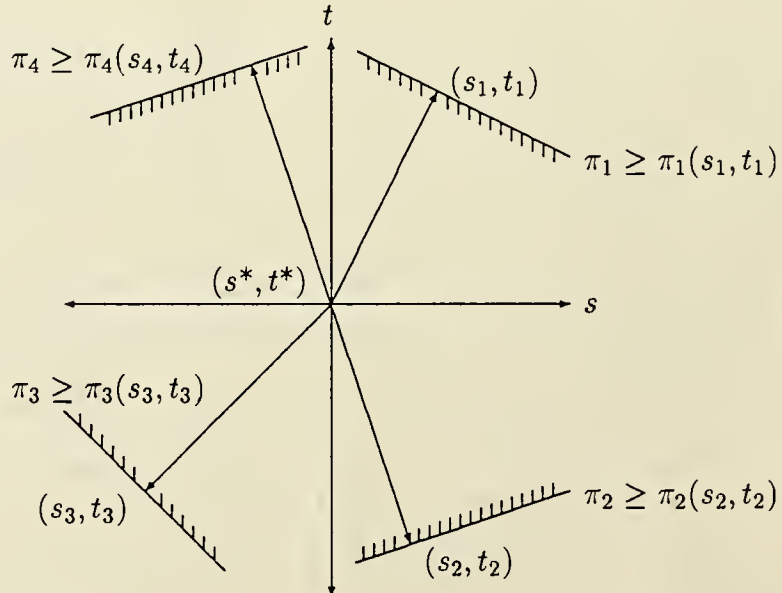
for $i \in \mathbf{N}$. Then (s_i, t_i) is also strictly individually rational for small enough $\varepsilon > 0$. Moreover, if $j \neq k$, then in the two-dimensional case,

$$\begin{aligned} \pi_j(s_k, t_k) &= \pi_j(s^*, t^*) - \varepsilon(\beta_j \beta_k + \gamma_j \gamma_k) / \sqrt{\beta_k^2 + \gamma_k^2} \\ &> \pi_j(s^*, t^*) - \varepsilon \sqrt{\beta_j^2 + \gamma_j^2} \\ &= \pi_j(s_j, t_j), \end{aligned}$$

where the inequality is due to Cauchy-Schwartz, and is strict by (P). Alternatively, in the negatively-sloped one-dimensional case,

$$\begin{aligned} \pi_j(s_k, t_k) &= \pi_j(s^*, t^*) - \varepsilon \theta_{jk} \sqrt{\beta_k^2 + \gamma_k^2} \\ &> \pi_j(s^*, t^*) - \varepsilon \sqrt{\beta_j^2 + \gamma_j^2} \\ &= \pi_j(s_j, t_j), \end{aligned}$$

since $\theta_{jk} < 0$ by (P). Finally, in the process it was established that $\pi_j(s_j, t_j) < \pi_j(s^*, t^*)$, as required. Illustrated below is the parameterization of the case of $n = 4$ players, with V_{24} a negatively-sloped line, and all other V_{ij} two-dimensional.



³This is nearly the converse of the first lemma in Abreu-Dutta (1991), who show that condition (iii) implies (P).

3. THE FOLK THEOREMS

We now turn to the folk theorems. That is, we show that any $u \in V^*$ can be (approximately) attained as a subgame perfect discounted average payoff vector in each of three distinctly different repeated contexts. Throughout, a^* is a correlated outcome yielding the desired payoff vector u .

Case A: Infinitely Repeated Games

Let $G(\delta)$ denote the infinitely-repeated game with stage game G and discount factor δ . Below is a modification of Theorem 2 in F-M, substituting (P) for full-dimensionality. The constructed equilibrium also differs from F-M in another sense, as it is *resilient*: After any finite sequence of deviations, play always returns to the equilibrium path within finitely-many periods.

3.1 INFINITE-HORIZON n -PLAYER FOLK THEOREM

Let $u = (u_1, \dots, u_n) \in V^*$. Suppose that (P) holds. Then $\exists \delta_0 < 1$ so that $\delta \in [\delta_0, 1] \implies G(\delta)$ has a subgame perfect discounted average payoff u .

Proof:

Along the equilibrium path, a^* is played. But if someone deviates then punishment interlude follows, consisting of a minimax phase and then a ‘recovery’ phase. Thus, the equilibrium strategies are: ⁴

1. Play a^* . [If player j deviates, start 2.]
2. Play M^j for Q periods. [If player $k \neq j$ deviates, start 3.] Then set $k \leftarrow j$.
3. Play p_k^* for R periods. [If any player j deviates, restart 2.] Return to step 1.

We next select Q and R so that this is a subgame perfect equilibrium. Let β and ω be the best and worst payoffs for any player in G , and first suppose that $\delta = 1$. For each step, we consider the ‘worst-case scenario’, where the incentive to deviate is greatest.

First note that given continuity of discounted sums in δ , if each deterrent is strict by some positive margin, say 1, they will remain strict for any level of discounting $\delta \in [\delta_0, 1]$, for some $\delta_0 < 1$. Now choose Q so that⁵

$$\omega + QU_j(p_j^*) > \beta + 1 \tag{1}$$

for all $j \in N$. Since $0 < U_j(p_j^*) < u_j$ by conditions (i) and (ii), (2) simultaneously renders the punishment interlude a strict deterrent to deviations from steps 1 and

⁴Throughout the paper, j , k , and l denote arbitrary players. Moreover, for clarity, we use the simple computer science $k \leftarrow j$ to mean “assign k the value j .” Also, program steps always follow sequentially, unless otherwise indicated. Conditional interrupts within square brackets are executed immediately.

⁵The inequality (1) — in particular, why the left-hand side is *not* $(Q + 1)U_j(p_j^*)$ — reflects the fact that obeying the random correlated outcome p_j^* might sometimes require j to play his worst possible outcome.

3, for any R . Next, step 3 deters deviations by the punishers from step 2 if R is large enough that

$$Q\omega + RU_k(p_j^*) > \beta + RU_k(p_k^*) + (Q - 1)u_k + 1 \quad (2)$$

for all $j, k \in \mathbf{N}$ with $j \neq k$. By condition (iii), such an R indeed exists. QED⁶

Case B: Finitely Repeated Games

In this section, we let $G(\delta, T)$ denote the T -fold repetition of G with the discount factor $\delta \leq 1$. What follows is a strengthening of Theorem 3.7 of B-K,⁷ on several fronts. First, we admit payoff discounting.⁸ Second, we substitute (P) for full-dimensionality. Finally, we weaken the condition that *all* players must have distinct Nash payoffs, to *all but one* — provided his Nash payoff is strictly positive. However, unlike B-K's (more intuitive) use of long deterministic cycles, we simply rely upon correlated outcomes. This is done to underscore the essential similarity among all three folk theorems. In fact, with finite-lived players, we essentially only need adjust the strategies of Case A to avoid the axe of backward induction. Finally, we remark that correlation also has the useful biproduct of permitting an *exact* (rather than approximate) folk theorem.

Let the Nash outcome e_i^* yield player $i \in \mathbf{N}$ the highest among all Nash payoffs of G , and f_i^* the lowest, for all $i \in \mathbf{N}$.

3.2 EXACT UNIFORM n -PLAYER FINITE-HORIZON FOLK THEOREM

Let $u = (u_1, \dots, u_n) \in V^*$. Suppose that (P) holds, and that either (a) $U_k(e_k^*) > U_k(f_k^*) \quad \forall k \in \mathbf{N}$, or (b) $U_k(e_k^*) > U_k(f_k^*) \quad \forall k \neq 1$, but $U_1(e_1^*) = U_1(f_1^*) > 0$. Then $\exists T_0 < \infty$ and $\delta_0 < 1$ so that $T \geq T_0$ and $\delta \in [\delta_0, 1] \implies G(\delta, T)$ has a subgame perfect discounted average payoff u .

Proof:

Let e^* be a correlated equilibrium according equal weight $1/n$ to each of the *preferred* Nash outcomes e_i^* . The typical T -period equilibrium outcome sequence is $\bar{a}, \dots, \bar{a}; e^*, \dots, e^*$, where e^* lasts S periods, and $U(\bar{a}) \in V^*$ will be seen to satisfy the required target payoff equation

$$(T - S)U(\bar{a}) + SU(e^*) = TU(a^*) \quad (3a)$$

if $\delta = 1$, or

$$(1 - \delta^{T-S})U(\bar{a}) + \delta^{T-S}(1 - \delta^S)U(e^*) = (1 - \delta^T)U(a^*) \quad (3b)$$

if $\delta < 1$. For ease of exposition, *late* deviations are those occurring during the final $Q + R + S$ periods of the repeated game; all others are called *early* deviations.

⁶Abreu-Dutta (1991) prove that (P) is also necessary for the infinite-horizon result.

⁷A much better and constructive demonstration is found in Krishna (1989): B-K requires a messy iterative approximation process.

⁸In so doing, we can establish a *uniform* folk theorem, meaning that the discount factor and horizon length can vary independently over the relevant range.

1. Play \tilde{a} until period $T - S$. [If player j deviates early, start 2; if player l deviates late, start 4.] Then play e^* until the end.
2. Play M^j for Q periods. [If player $k \neq j$ deviates, start 3.] Then set $k \leftarrow j$.
3. Play p_k^* for R periods. [If some player j deviates early, restart 2; if player $l \neq 1$ deviates late, start 4. If 1 deviates late, start 5.] Then return to step 1.
4. Play f_l^* until the end.
5. Play M^1 until period $T - S + \lfloor \sqrt{S} \rfloor$. [If player $l \neq 1$ deviates, start 4.] Then play e^* until the end.

We proceed recursively. First choose Q , R , and S to ensure subgame perfection. ‘Early’ deviations are handled exactly as in the infinite horizon case. Thus (1) determines Q , while (2) must be modified because $\tilde{a} \neq a^*$ (unless perchance $a^* = e^*$). We now require that

$$Q\omega + RU_k(p_j^*) > \beta + RU_k(p_k^*) + (Q - 1)U_k(\tilde{a}) + 1 \quad (4)$$

for all $k \neq j$. Given that \tilde{a} is as yet unspecified, we shall instead insist that R satisfy

$$Q\omega + RU_k(p_j^*) > \beta + RU_k(p_k^*) + (Q - 1)(2u_k - U_k(p_k^*)) + 1. \quad (5)$$

This will turn out to imply (4). Moreover, \tilde{a} will be chosen so that the personalized punishment vector p_k^* satisfies the analogue of condition (ii) with respect to \tilde{a} , $\forall k \in \mathbf{N}$. Thus, just as in Case A, the punishment interlude, steps 2 and 3, will deter deviations from step 1.

Next, step 4 will deter all ‘late’ deviations by players $l \neq 1$ so long as S is large enough that

$$(Q + R + \sqrt{S})\omega + (S - \sqrt{S})U_j(e^*) > \beta + (Q + R + S - 1)U_j(f_j^*) + 1 \quad (6a)$$

for all $l \neq 1$. This inequality obtains for all sufficiently large S , because each player $l \neq 1$ strictly prefers e^* to f_l^* , and since S grows faster than \sqrt{S} . Similarly, step 5 is a deterrent to late deviations by player 1 provided

$$Q\omega + RU_1(p_1^*) + SU_1(e^*) > \beta + (S - \sqrt{S})U_1(e^*) + 1 \quad (6b)$$

and

$$SU_1(e^*) > \beta + (S - \sqrt{S})U_1(e^*) + 1. \quad (6c)$$

Clearly, for large enough S , inequalities (6a), (6b), and (6c) are valid.

Given Q , R , and S as defined above, the above program is feasible if $T_0 \geq Q + R + S$. Next, since $a^* \in V^*$, there is some η -neighbourhood around u entirely contained within V^* . Then let T_0 be sufficiently large that

$$\frac{S}{T_0 - S} \|u - U(e^*)\| < \min\langle \eta, \min_{j \in \mathbf{N}} [u_j - U_j(p_j^*)] \rangle \quad (7a)$$

so that

$$\delta^{T-s} \frac{1 - \delta^S}{1 - \delta^{T-S}} \|u - U(e^*)\| < \min\langle \eta, \min_{j \in \mathbf{N}} [u_j - U_j(p_j^*)] \rangle \quad (7b)$$

for all $\delta \leq 1$ and $T \geq T_0$.

Finally, let $T \geq T_0$, and introduce discounting just as in Case A, so that each deterrent remains strict for all $\delta \in [\delta_0, 1]$, for some $\delta_0 < 1$. Define \tilde{a} implicitly by the target payoff equation (3a) or (3b). To verify that indeed $U(\tilde{a}) \in V^*$, (3a) and (3b) can be rewritten as

$$U(\tilde{a}) - u = \begin{cases} \frac{S}{T-S} [u - U(e^*)] & \text{if } \delta = 1 \\ \delta^{T-s} \frac{1 - \delta^S}{1 - \delta^{T-S}} [u - U(e^*)] & \text{if } \delta < 1 \end{cases} \quad (8)$$

In light of (7a) and (7b), equation (8) both says that $\tilde{a} \in V^*$ and also that

$$|U_k(\tilde{a}) - u_k| < \min_{j \in \mathbf{N}} [u_j - U_j(p_j^*)] \quad (9)$$

for all $k \in \mathbf{N}$. By the triangle inequality, there are two immediate consequences of (9). First, we may conclude that $U_k(\tilde{a}) < 2u_k - U_k(p_k^*)$, so that (5) implies (4); and second, that

$$\begin{aligned} U_k(\tilde{a}) - U_k(p_k^*) &= u_k - U_k(p_k^*) + U_k(\tilde{a}) - u_k \\ &\geq [u_k - U_k(p_k^*)] - |U_k(\tilde{a}) - u_k| > 0, \end{aligned}$$

both of which were asserted earlier.

QED

Case C: Overlapping Generations Games

Smith (1990) considers a model of overlapping generations games $\text{OLG}(G; \delta, T)$, in which — in its most basic formulation — a player dies and is replaced every T periods.⁹ Using the techniques developed for Case B, it is possible to render that folk theorem *exact* too. Moreover, its full-dimensionality condition can also be weakened to (P).

3.3 EXACT UNIFORM n -PLAYER OLG FOLK THEOREM

Let $u = (u_1, \dots, u_n) \in V^*$. Suppose that (P) holds. Then $\exists T_0 < \infty$ and $\delta_0 < 1$ so that $T \geq T_0$ and $\delta \in [\delta_0, 1] \implies \text{OLG}(G; \delta, T)$ has a subgame perfect discounted average payoff u .

A parallel proof of this result is possible, but is omitted.

We remark that (P) is required here — as elsewhere — to payoff distinguish between players *at the stage game level*: In Cases A and B, this was the only tool at our disposal. But in overlapping generations games, because the players' tenures in the game do not coincide, we may also distinguish between them *intertemporally*. That is, we may await the death of specific players before rewarding or punishing

⁹Kandori (1990) describes a related model, but doesn't consider the uniform folk theorem discussed below.

the others.¹⁰ In the as yet unexplored arena of *spatial* games, one might distinguish among the players it interspatially. This insight allows us to more generally think of (P) as just one method of awarding players separable payoff streams.

4. CONCLUDING REMARKS

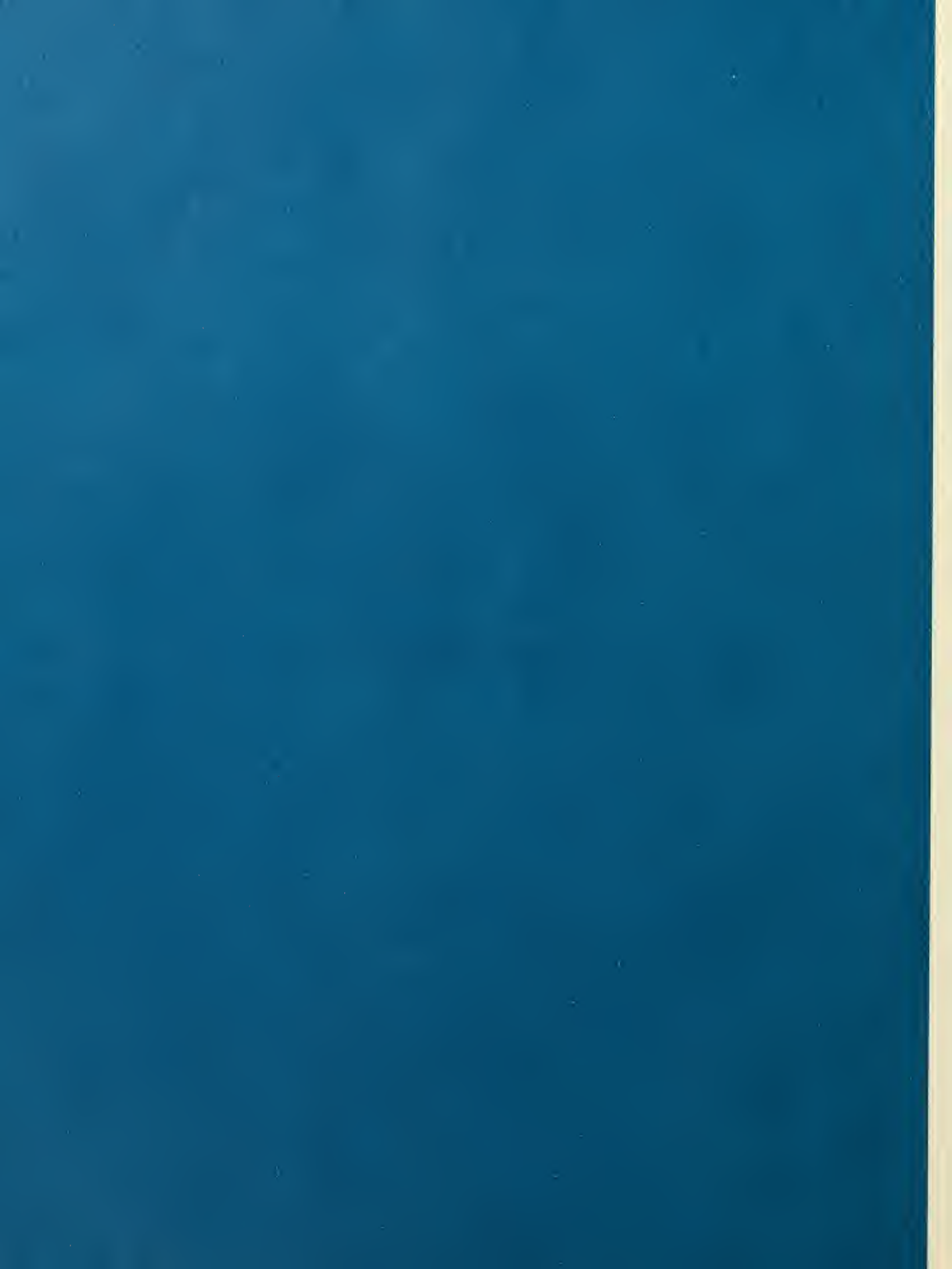
The requirement that games with $n > 2$ players have a full-dimensional payoff space is unnecessarily harsh. Although it must be noted that full-dimensionality has an obvious interpretation itself as an assumption of genericity, our new condition (P) can be easily seen to hold generically for the smaller class of (stage) games for whose payoff spaces have dimensions $2, 3, \dots, n - 1$. Such a class might prove economically relevant — for instance, when the payoffs of all players or even of several coalitions has constant sum.

Not only is our condition (P) less stringent, but it is also *more natural*. For it emphasizes the fundamental basis of all Nash strategies, that it is only necessary to payoff distinguish *any two* players at once, and not all of them. Deviants can be singled out one at a time, so that it is pure overkill to insist that all punishers be rewarded in the stage game. Such a requirement ignores the wealth of available *dynamic* strategies in a repeated game.

REFERENCES

- Abreu, D. and P. Dutta, A Folk Theorem for Discounted Repeated Games: A New Condition, Mimeo (University of Rochester, Rochester, NY).
- Benoit J.-P. and V. Krishna, 1985, Finitely Repeated Games, *Econometrica* 53, 905-922.
- Fudenberg, D. and E. Maskin, 1986, The Folk Theorem in Repeated Games with Discounting or with Incomplete Information, *Econometrica* 54, 533-554.
- Kandori, M., 1990, Repeated Games Played by Overlapping Generations of Players, forthcoming in *Review of Economic Studies*.
- Krishna, 1989, The Folk Theorems for Repeated Games, Mimeo 89-003, Harvard Business School.
- Rubinstein, A., 1979, Equilibrium in Supergames with the Overtaking Criterion, *Journal of Economic Theory* 21, 1-9.
- Smith, L., 1990, Folk Theorems in Overlapping Generations Games, forthcoming in *Games and Economic Behavior*.

¹⁰The development of this idea is the substance of the *non-uniform* folk theorems of Smith (1990) and Kandori (1990).



Date Due 9-8-92

MAY 13 1992

MIT LIBRARIES DUPL



3 9080 00756965 7

