

VALUE AND PERFECTION IN STOCHASTIC GAMES

BY

ROBERT SAMUEL SIMON*

*Faculty of Mathematics, London School of Economics**Houghton Street, London WC2A 2AE, U.K.**e-mail: R.S.Simon@lse.ac.uk*

ABSTRACT

A stochastic game is **valued** if for every player k there is a function $r^k: S \rightarrow \mathbf{R}$ from the state space S to the real numbers such that for every $\epsilon > 0$ there is an ϵ equilibrium such that with probability at least $1 - \epsilon$ no state s is reached where the future expected payoff for any player k differs from $r^k(s)$ by more than ϵ . We call a stochastic game **normal** if the state space is at most countable, there are finitely many players, at every state every player has only finitely many actions, and the payoffs are uniformly bounded and Borel measurable as functions on the histories of play. We demonstrate an example of a recursive two-person non-zero-sum normal stochastic game with only three non-absorbing states and limit average payoffs that is not valued (but does have ϵ equilibria for every positive ϵ). In this respect two-person non-zero-sum stochastic games are very different from their zero-sum varieties. N. Vieille proved that all such non-zero-sum games with finitely many states have an ϵ equilibrium for every positive ϵ , and our example shows that any proof of this result must be qualitatively different from the existence proofs for zero-sum games. To show that our example is not valued we need that the existence of ϵ equilibria for all positive ϵ implies a “perfection” property. Should there exist a normal stochastic game without an ϵ equilibrium for some $\epsilon > 0$, this perfection property may be useful for demonstrating this fact. Furthermore, our example sows some doubt concerning the existence of ϵ equilibria for two-person non-zero-sum recursive normal stochastic games with countably many states.

* This research was supported financially by the German Science Foundation (Deutsche Forschungsgemeinschaft) and the Center for High Performance Computing (Technical University, Dresden). The author thanks Ulrich Krengel and Heinrich Hering for their support of his habilitation at the University of Goettingen, of which this paper is a part.

Received December 23, 2004 and in revised form July 21, 2005

1. Introduction

An **equilibrium** is a set of strategies, one for each player, such that no player can gain in payoff by choosing a different strategy, given that all the other players do not change their strategies. A two-person game is zero-sum when the payoff to one player is always the negation of the payoff to the other player; if there is an equilibrium for a zero-sum game then there is a unique equilibrium payoff (to Player One) and it is called the **value** of the game.

For any $\epsilon \geq 0$, an ϵ -**equilibrium** in a game is a set of strategies, one for each player, such that no player can gain in payoff by more than ϵ by choosing a different strategy, given that all the other players do not change their strategies. We say that approximate equilibria exist if for every $\epsilon > 0$ there exists an ϵ -equilibrium. If a game has no equilibrium but does have approximate equilibria, then it is unavoidable that there will be some advantage to some player to break from the prescribed behavior; however, this advantage can be made as small as one wants. An equilibrium is often too much to expect, but approximate equilibria are the next best thing. Given finitely many players and a uniform bound on the payoffs, any vector cluster point of ϵ equilibrium payoffs as ϵ converges to zero will give to the game a kind of equilibrium payoff to each player. If the game is also zero-sum and has approximate equilibria, then there will be a unique cluster point of ϵ equilibrium payoffs (for Player One) as ϵ converges to zero, also called the value of the game.

A **stochastic** game is played on a state space. The present state and the present behavior of all players determines stochastically the transition to a new state. All players have complete knowledge of the past history of play. A priori there is no bound on the number of stages of play.

We define a stochastic game to be **normal** if

- (1) there are countably many states,
- (2) there are finitely many players and at any state the action sets for all players are finite,
- (3) all the payoffs defined in the game are uniformly bounded,
- (4) the payoffs are functions on the histories of play that are measurable with respect to the Borel σ -algebra defined by the finite stages of the game.

This fourth property will be made more precise later.

This paper concerns only normal stochastic games. We are interested in the complexity generated endogenously between the definition of the game and its equilibria.

Shapley [7] introduced the concept of a stochastic game in the context of zero-

sum games where the payoffs are evaluated according to a function determined only by the state and the pair of actions taken and a discount factor. A discount factor is a positive quantity ρ strictly less than one such that for every stage $i \geq 0$ of play the $i + 1$ st stage is worth only ρ times that of the i th stage of play. The finiteness of the geometric sum $1 + \rho + \rho^2 + \dots$ gives the game a compact structure. Shapley showed that such discounted normal zero-sum games with finitely many states have equilibria and values obtainable from stationary strategies, meaning that the strategies are independent of both the stage and history of play (and are dependent only on the state).

Another way to evaluate the payoffs is by some limit of average values determined by the states and the actions chosen, averaged over the stages of play. This is called the **limit average** evaluation. For example, if a player received a sequence w_0, w_1, \dots of payoffs on the stages $0, 1, \dots$, respectively, then her payoff could be $\lim_{i \rightarrow \infty} \sup \frac{1}{i+1} \sum_{k=0}^i w_k$. In general, when payoffs are limit average normal stochastic games do not have equilibria. This was demonstrated by Blackwell and Ferguson [1] with their famous zero-sum example “The Big Match”. However, they showed that this game does have approximate equilibria and a value.

Mertens and Neyman [5] proved that every zero-sum normal stochastic game with limit average payoffs played on a finite state space has approximate equilibria and a value. Maitra and Sudderth [3] extended this result to countably many states and Martin [4] extended this result further to payoff functions defined on the infinite paths of play that are Borel with respect to their finite stage truncations.

Concerning two-player non-zero-sum games with limit average payoffs, the central result was accomplished by Vieille [8]; he proved that all such normal stochastic games with finitely many states have approximate equilibria. For two-player non-zero-sum normal stochastic games with countably many states the question is still open.

If the time horizon of a normal stochastic game is truncated so that there is a finite maximum number of stages then equilibria will exist [6]. The open-ended nature of the time horizon gives stochastic games their most important theoretical complication.

The property of “valued” for stochastic games (stated in the abstract) is a natural strengthening of the approximate equilibria property. Consider a sequence of ϵ_i equilibria with the ϵ_i converging to zero and let s be a state such that with these approximate equilibria s is reached with probabilities that do

not converge to zero. For every player k there will be a quantity $r^k(s)$ and a subsequence of the ϵ_i equilibria such that, conditioned on the first visit to s , the expected payoffs for player k will converge to $r^k(s)$. At the second visit to the state s the sub-game that remains is no different from the game that starts at this state. Of course the equilibria are only approximate, meaning that there could be a very small probability that the players find themselves in a situation at some state s calling for behavior that is far from that of an equilibrium for the game that starts at s . But one could expect some property of payoff stability for some sequence of ϵ equilibria like that of the valued property, namely that if a state is reached with large probability then almost all visits to this state yield approximately the same payoffs.

Before going further, an additional concept is necessary. A state is **absorbing** if once this state is reached the play can never leave this state, no matter what the players do. If zero is the payoff for all players at all non-absorbing states then the game is **recursive**. An absorbing state defines an isolated sub-game, and therefore in general one associates to an absorbing state a fixed payoff for each player representing what they receive in equilibrium if this absorbing state is reached. An absorbing state in a stochastic game demonstrates vividly the difference between finitely repeated games and discounted games on the one hand and limit average games on the other. In a finitely repeated game or a discounted game when an absorbing state is reached the payoff to a player would be a convex combination of the absorbing state payoffs and the payoffs on the stages before this absorbing state is reached. If the payoffs to the players are limit average, once reaching an absorbing state the payoffs to the players are determined by this absorbing state only.

In general normal recursive stochastic games do not possess the valued property. With three players and finitely many states Flesch, Thuijsman and Vrieze [2] found a counter-example. Their game is very simple, involving only one special state where the players have any influence on the outcome of the game. They showed that for sufficiently small $\epsilon > 0$ the only ϵ equilibria involve cyclic behavior, both of the players and of the payoffs conditioned on the event that the game has remained at the special state.

We introduce a recursive two-person non-zero-sum normal stochastic game with three non-absorbing states and limit average payoffs that is not valued. New about our example is that the lack of the valued property is possible with only two players (and finitely many states and actions).

In Vieille's proof of approximate equilibria (for two-player non-zero-sum

games with finitely many states) the behavior of the players is very complicated, involving nothing like the valued property. Naturally one wonders if there is a much simpler proof for Vieille's result, a proof closer in style to those for zero-sum normal stochastic games which do utilize the valued property, for example with some kind of fixed point argument on the space of payoff vectors (see [3]). Zero-sum normal stochastic games are valued; the proof is easy and provided below. Our main result shows that such an alternative approach is not possible; there is an aspect of Vieille's proof (or any alternative proof of Vieille's result) that must be complicated.

With the Vieille proof, the usual behavior of the players depends on more than the state which is visited. For some non-absorbing state that could be visited infinitely many times the players may return often to this state in the ϵ equilibria such that the expected number of visits is in inverse proportion to a fixed power of ϵ . At such a state signals may be given by one of the players, and the future behavior of both players may be dependent on these signals. Of course if the signals distinguish between two different ways to play such that in the limit both signals are given with large probability and they imply significantly different expected future payoffs to the player who is not giving the signal, then the valued property would be contradicted.

Zero-sum games have a monotonicity relationship between the strategies and their values. If the payoff to a player associated with a combination of actions in a zero-sum game is increased, then the value for that player cannot go down. This is very different with non-zero-sum games and their equilibria. By increasing the payoff to some player associated with a combination of actions, an equilibrium based on cooperation may be destroyed. The introduced lack of trust and the re-establishment of balance in all new equilibria could result in lower payoffs for the player whose payoff was increased. To obtain our main result, we exploit this lack of monotonicity.

In the next section we describe the model of normal stochastic games. We establish relationships between approximate equilibria, value, and perfection. The third section contains our example and the proof that it is not valued. The fourth and last section discusses the open problem of approximate equilibria for games with countably many states.

2. The model

We define normal stochastic games similarly to Markov chains and their harmonic functions. The additional complication concerns the freedom of the play-

ers to influence the transitions based on the previous history of play. By extending the definition of the state space so that distinct past histories of play lead to distinct states, we would stay close to the explicit context of Markov chains, as then a strategy choice from each player would define a Markov chain. But then we could lose track of the original structure of the game, especially if it is defined by a finite state space.

For every finite or countable set A let $\Delta(A)$ stand for the set of all probability distributions on A . If A is finite then $\Delta(A)$ is a finite-dimensional simplex. If $x \in \Delta(A)$ and $a \in A$ then the a coordinate of x will be represented as $x(a)$ (the probability given to a by x).

There is a countable or finite state space S and a finite set N of players. For every player $n \in N$ and every $s \in S$ there is a finite set A_s^n of actions. (If the action sets are countable then even for one-stage zero-sum games there are examples without approximate equilibria, for example the game where the two players choose natural numbers and the player who chooses the larger number wins a unit value.) For every $s \in S$ and every $a \in A_s := \prod_{n \in N} A_s^n$ (a choice of action for each player) there will be a transition law $p_a^s \in \Delta(S)$ governing the transition to states at the next stage of play after a visit to s .

We assume that the game starts at an initial state $\hat{s} \in S$. (If one prefers to start with a distribution on all the states in S one can add an initial state \hat{s} that occurs only at the start of the game and such that every player has only one action at this state.) Define

$$\mathcal{H}_\infty := \{(\hat{s} = s_0, a_0, s_1, a_1, \dots) \mid \forall i \geq 0 a_i \in A_{s_i}, p_{a_i}^{s_i}(s_{i+1}) > 0\},$$

the set of infinite histories of play. For the initial state $\hat{s} \in S$ let $\mathcal{H}_0^{\hat{s}} := \{(\hat{s})\}$, and for every $i \geq 1$ let \mathcal{H}_i^s be the set of truncations of \mathcal{H}_∞ of the form $(\hat{s} = s_0, a_0, s_1, a_1, \dots, s_{i-1}, a_{i-1}, s_i = s)$ (leaving out the actions at stage i). Let \mathcal{H}_i be the union $\bigcup_{s \in S} \mathcal{H}_i^s$ and let \mathcal{H}^s be the union $\bigcup_{i=0}^\infty \mathcal{H}_i^s$ (with \mathcal{H}_i^s the empty set if s is not reachable on stage i). Let \mathcal{H}_ω be the union $\bigcup_{i=0}^\infty \mathcal{H}_i = \bigcup_{s \in S} \mathcal{H}^s$. If $h \in \mathcal{H}_\omega$ is also in \mathcal{H}^s then we say that h **terminates** at s .

A payoff for a player $n \in N$ in a normal stochastic game is a function \mathcal{V}^n on \mathcal{H}_∞ that is uniformly bounded and measurable with respect to the Borel σ -algebra generated by the partitions on \mathcal{H}_∞ induced by the \mathcal{H}_i . A two-player game is **zero-sum** if $\mathcal{V}^1(h) + \mathcal{V}^2(h) = 0$ for all $h \in \mathcal{H}_\infty$ (where without loss of generality we assume that $N = \{1, 2\}$). Let $M > 1$ be a positive real number larger than the maximal difference between all payoffs in the game.

A strategy σ^n of Player $n \in N$ is a collection of functions $(\sigma_s^n \mid s \in S)$ such that for every $s \in S$, σ_s^n is a function from \mathcal{H}^s to $\Delta(A_s^n)$. For every

tuple of strategies $\sigma = (\sigma^n \mid n \in N)$, one strategy for each player, probability distributions $\mu_{\sigma,i}$ are induced on the \mathcal{H}_i in the natural way. We start at the initial history $(\hat{s}) \in \mathcal{H}_0^{\hat{s}}$ with $\mu_{\sigma,0}(\{(\hat{s})\}) = 1$. Given that $\mu_{\sigma,i}(h_i)$ is positive for some $h_i \in \mathcal{H}_i^{s_i}$ and $h_{i+1} \in \mathcal{H}_{i+1}$ is a history such that the i stage truncation of h_{i+1} is equal to $h_i \in \mathcal{H}_i^{s_i}$ with $h_{i+1} = (h_i, a_i, s_{i+1})$ and $a_i = (a_i^n \mid n \in N)$, we define inductively $\mu_{\sigma,i+1}(h_{i+1}) := \mu_{\sigma,i}(h_i) p_{a_i}^{s_i}(s_{i+1}) \prod_{n \in N} \sigma_{s_i}^n(h_i)(a_i^n)$. A Borel probability distribution μ_σ is induced on \mathcal{H}_∞ in the natural way, by the $\mu_{\sigma,i}$ and Kolmogorov's Extension Theorem. For every player $n \in N$ and every strategy tuple σ the distribution μ_σ generates a payoff $\mathcal{V}^n(\sigma)$ for player n as the expected value of the function \mathcal{V}^n on \mathcal{H}_∞ , determined by the probability distribution μ_σ .

For any tuple $\sigma = (\sigma^n \mid n \in N)$ of strategies, an alternative tuple $\tilde{\sigma} = (\tilde{\sigma}^n \mid n \in N)$ and a player $k \in N$, define $\sigma|\tilde{\sigma}^k$ to be the tuple such that $\tilde{\sigma}^k$ is the strategy for player k but if $n \neq k$ then σ^n is the strategy for player n . An ϵ equilibrium is a strategy tuple $\sigma = (\sigma^n \mid n \in N)$ such that for any alternative tuple $(\tilde{\sigma}^n \mid n \in N)$ and every player $n \in N$ it holds that $\mathcal{V}^n(\sigma|\tilde{\sigma}^n) \leq \epsilon + \mathcal{V}^n(\sigma)$. A zero-sum game has the value $r \in \mathbf{R}$ for a designated first player if for every positive ϵ there is an ϵ equilibrium whose expected payoff for the first player is within ϵ of r .

A stochastic game is a **limit average** game when for every player $n \in N$ the Borel function \mathcal{V}^n is between $\lim_{i \rightarrow \infty} \inf$ and $\lim_{i \rightarrow \infty} \sup$ of the average $\frac{1}{i} \sum_{k=0}^{i-1} w_{s_k}^n(a_k)$ where, for every state $s \in S$ and $n \in N$, w_s^n is a real function defined on $A_s = \prod_{m \in N} A_s^m$.

For any tuple σ of strategies, a player $n \in N$, and a stage i of play define $v_\sigma^n: \mathcal{H}_i \rightarrow \mathbf{R}$ by $v_\sigma^n(h_i)$ equaling the expected value of $\mathcal{V}^n(\sigma)$ conditioned on reaching h_i on the i th stage, with $v_\sigma^n(h_i)$ defined to be any quantity bounded within the payoffs defining the game if the probability of reaching h_i is zero. Extend the definition of v_σ^n to \mathcal{H}_ω . For any fixed σ and player $n \in N$ the function v_σ^n on \mathcal{H}_ω with respect to μ_σ is harmonic.

Define a stochastic game to be **valued** if for every player $n \in N$ there exists a function $r^n: S \rightarrow \mathbf{R}$ such that for every $\epsilon > 0$ there is an ϵ equilibrium σ to the game such that the probability does not exceed ϵ that some history $h_i \in \mathcal{H}_i^s$ occurs with $|v_\sigma^n(h_i) - r^n(s)| > \epsilon$ for some player $n \in N$.

PROPOSITION 1: *All zero-sum normal stochastic games are valued.*

Proof: Let $r: S \rightarrow \mathbf{R}$ be defined so that $r(s)$ is the value of the game (for Player One) starting at the state s . (The existence of this value was proven by Martin [4].) For any $0 < \epsilon \leq 1$ let σ be a strategy pair where both players

guarantee a payoff (for Player One) within $\epsilon^2/10$ of the value $r(\hat{s})$, where \hat{s} is the initial state. Assuming that $\max_i |v_\sigma^1(h_i) - r(s_i)| > \epsilon$ is obtained for a subset $\mathcal{A} \subseteq \mathcal{H}_\omega$ of finite histories reached with a probability of more than ϵ (according to the distribution μ_σ), we can assume by symmetry that there is a subset $\mathcal{A}_1 \subseteq \mathcal{A}$ given a probability of at least $\epsilon/2$ such that $\min_i v_\sigma^1(h_i) - r(s_i) < \epsilon$ is obtained (since otherwise we obtain the same for Player Two with $-r$ replacing r). We alter the strategy of Player One so that at the first stage i such that $v_\sigma^1(h_i) - r(s_i) < \epsilon$ Player One switches her strategy to one guaranteeing $r(s_i) - \epsilon/2$. The gain in expected payoff for Player One would be at least $\epsilon^2/4$, a contradiction. ■

For every player n define $\chi^n: S \rightarrow \mathbf{R}$ so that $\chi^n(s)$ is the min-max value for player n at the state s , the upper bound for what player n can obtain from a start at s in response to all strategy choices of the other players. Formally $\chi^n(s)$ equals $\inf_\sigma \sup_{\tilde{\sigma}^n} \mathcal{V}_s^n(\sigma|\tilde{\sigma}^n)$, where the payoff function \mathcal{V}_s^n is defined by the game for which s is the initial state. The importance of the min-max value χ^n is that it represents the ability of the players to punish player n with pre-determined strategies (for example as part of an approximate equilibrium). Because the other players may be limited in their ability to coordinate their actions, this min-max value could be strictly greater than the max-min value when there are at least three players. (The max-min value for player n is the most that player n can obtain when she must choose her strategy first and the other players respond to that strategy choice with the goal of minimizing her payoff, obtained formally by switching the inf and the sup in the above formula.) Proofs that (two-person) zero-sum games have approximate equilibria demonstrate that for these games evaluated at the initial state the max-min value equals the min-max value, thereafter called the value of the game.

For every $a^n \in A_s^n$ and $\hat{a} \in \prod_{k \neq n} A_s^k$ let (\hat{a}, a^n) be the corresponding member of $A_s = \prod_{k \in N} A_s^k$, with \hat{a}^k the corresponding action of Player k for all $k \neq n$. For every player $n \in N$ and strategy tuple σ , define the **jump** function $j_\sigma^n: \mathcal{H}_\omega \rightarrow \mathbf{R}$ by

$$j_\sigma^n(h) = \max_{a^n \in A_s^n} \sum_{t \in S} \chi^n(t) \sum_{\hat{a} \in \prod_{k \neq n} A_s^k} \prod_{k \in N \setminus \{n\}} \sigma_s^k(h)(\hat{a}^k) p_{(\hat{a}, a^n)}^s(t),$$

namely the maximal expected value of χ^n on the next stage following s . Extend this definition to $j_\sigma^n: \mathcal{H}_\omega \rightarrow \mathbf{R}$ in the natural way.

For any function $f: \mathcal{H}_\omega \rightarrow \mathbf{R}$, state $s \in S$, finite history $h \in \mathcal{H}^s$, action $a^n \in A_s^n$ and strategy tuple σ , define $w_\sigma^f(h)(a^n)$ to be the expected value of f on the next stage after h , conditioned on the use of a^n by Player j and the use

of $\sigma_s^k(h)$ by all the other players $k \neq j$. This means that

$$w_\sigma^f(h)(a^n) = \sum_{t \in S} \sum_{\hat{a} \in \prod_{k \neq n} A_s^k} f(h, (\hat{a}, a^n), t) \prod_{k \in N \setminus \{n\}} \sigma_s^k(h)(\hat{a}^k) p_{(\hat{a}, a^n)}^s(t).$$

Define $w_\sigma^n(h)(\hat{a}^n)$ to be $w_\sigma^{v_\sigma^n}(h)(a^n)$.

Definition: A strategy tuple σ of a stochastic game is ϵ **perfect** if for every player $n \in N$ there exists a function $r^n: \mathcal{H}_\omega \rightarrow \mathbf{R}$ and a subset $\mathcal{B} \subseteq \mathcal{H}_\omega$ such that the probability of reaching $\mathcal{H}_\omega \setminus \mathcal{B}$ with the strategies σ does not exceed ϵ and for all players $n \in N$ and all finite histories $h \in \mathcal{B}$,

$$r^n(h) \geq j_\sigma^n(h) - \epsilon, \\ |r^n(h) - v_\sigma^n(h)| \leq \epsilon, \quad \text{and}$$

for all actions a^n chosen with positive probability by σ^n at h

$$|w_\sigma^n(h)(a^n) - r^n(h)| \leq \epsilon.$$

Furthermore, the strategy tuple σ is ϵ **value-perfect** if for every player n the function r^n is dependent only on the state where the history terminates. A stochastic game is **perfect** if there exists an ϵ perfect strategy tuple for every positive ϵ . A stochastic game is **value-perfect** if there exists a ϵ value-perfect strategy tuple for every positive ϵ .

THEOREM 1: *A normal stochastic game with approximate equilibria is also perfect. A valued normal stochastic game is also value-perfect.*

Proof: For a fixed $0 < \epsilon \leq 1$ let σ be an $\epsilon^4/(4000|N|^4M^2)$ equilibrium, and if the game is valued let $r^n: S \rightarrow \mathbf{R}$ be the value functions and let \mathcal{A} be a subset of finite histories such that \mathcal{A} is reached with a probability of no more than $\epsilon^4/(4000|N|^4M^2)$ and such that if h is in $\mathcal{H}_\omega \setminus \mathcal{A}$ terminating at $s \in S$ then $|v_\sigma^n(h) - r^n(s)| \leq \epsilon^4/(4000|N|^4M^2)$ for all $n \in N$. Otherwise let \mathcal{A} be the empty set.

For any history $h \in \mathcal{H}_i^s$ and player n let $A_+^n(h)$ be the set of actions $a^n \in A_s^n$ chosen with positive probability according to $\sigma_s^n(h)$ such that $w_\sigma^n(h)(a^n) - v_\sigma^n(h) > \epsilon/10$ and let $A_-^n(h)$ be the set of actions $a^n \in A_s^n$ chosen with positive probability according to $\sigma_s^n(h)$ such that $w_\sigma^n(h)(a^n) - v_\sigma^n(h) < -\epsilon/10$. Let $A^n(h)$ be the union of $A_+^n(h)$ with $A_-^n(h)$. Whenever $A^n(h)$ is not empty player n could alter her strategy in the following way. If a^n is in $A_+^n(h)$ then a^n could be chosen with certainty (or any other action in this set). If a^n is in $A_-^n(h)$ then

the probability given to a^n could be given to any action \hat{a}^n maximizing the value of $w_\sigma^n(h)(\cdot)$. In either case we conclude that player n can increase her expected payoff conditioned on reaching h by at least $\epsilon\sigma_s^n(h)(a^n)/10$ (this much if a^n is in $A_-^n(h)$ and by at least $\epsilon/10$ if a^n is in $A_+^n(h)$). Because σ is an $\epsilon^4/(4000|N|^4M^2)$ equilibrium, the probability of any player n ever using an action in some $A^n(h)$ according to the distribution μ_σ does not exceed $\epsilon^3/(400|N|^3M^2)$.

Let \mathcal{D}_1 be the union of \mathcal{A} with the histories h in \mathcal{H}_ω , where for some player $n \in N$ an action in $A^n(h_i)$ is used at some history h_i resulting from the i th stage truncation of h (where $h \in \mathcal{H}_k$ and $i \leq k$). By the above we know that the probability of reaching \mathcal{D}_1 according to μ_σ does not exceed $\epsilon^3/(350|N|^3M^2)$.

For any history $h \in \mathcal{H}_i^s$ and player $n \in N$ let $B^n(h)$ be the set of actions $a^n \in A_s^n$ chosen with positive probability according to σ_s^n such that the probability of entering \mathcal{D}_1 on any following stage is at least $\epsilon/(10|N|M)$ when a^n is used against the distributions $\sigma_s^k(h)$ for $k \neq n$ and otherwise at later stages all players behave according to σ . $A^n(h)$ is a subset of $B^n(h)$ (since use of an action in $A^n(h)$ causes entry into \mathcal{D}_1 with certainty). By the above the probability of any player n ever using an action in some $B^n(h)$ according to μ_σ does not exceed $\epsilon^2/(35M|N|^2)$.

We define a new strategy tuple $\hat{\sigma}$. If $h \in \mathcal{H}_i^s$ and no member of $B^n(h')$ appears in h for any truncation $h' \in \mathcal{H}_k$ of h with $k < i$, then $\hat{\sigma}_s^n(h)$ is the distribution on A_s^n where the actions in $B^n(h)$ are given zero probability and the probabilities for the remaining actions are normalized (meaning that the probability for an action $a^n \notin B^n(h)$ is $\sigma_s^n(h)(a^n)/\sum_{\bar{a}^n \notin B^n(h)} \sigma_s^n(h)(\bar{a}^n)$; if all actions are removed in this way, then define $\hat{\sigma}_s^n(h)$ to be any distribution). Otherwise, if $B^n(h)$ is empty or some member of $B^n(h')$ was played in the past then $\hat{\sigma}_s^n(h) = \sigma_s^n(h)$.

Let \mathcal{D}_2 be the subset of \mathcal{H}_ω such that $h \in \mathcal{D}_2$ if and only if there is some truncation h_k of h and some player n with $|v_\sigma^n(h_k) - v_{\hat{\sigma}}^n(h_k)| > \epsilon/5$. From the unlikeliness of using an action in some $B^n(h)$ we conclude that the probability of reaching \mathcal{D}_2 according to μ_σ does not exceed $\epsilon/7$.

Let \mathcal{D}_3 be the subset of \mathcal{H}_ω defined by $h \in \mathcal{D}_3$ if and only if for some player n and some truncation $h_k \in \mathcal{H}_k$ of h the actions removed to make $\hat{\sigma}^n(h_k)$ from $\sigma^n(h_k)$ had a probability greater than $\epsilon/(5|N|M)$. Also from the unlikeliness of using an action in some $B^n(h)$ the probability of reaching \mathcal{D}_3 according to μ_σ does not exceed $\epsilon/7$.

Let \mathcal{D}_4 be the subset of \mathcal{H}_ω defined by $h \in \mathcal{D}_4$ if and only if there is some player n and some truncation $h_k \in \mathcal{H}_k$ of h such that $v_\sigma^n(h_k) < j_\sigma^n(h_k) - \epsilon/10$.

Whenever this inequality holds player n can obtain an expected payoff of at least $v_\sigma^n(h_k) + \epsilon/11$, and therefore by the approximate equilibrium property of σ the probability of reaching \mathcal{D}_4 according to μ_σ cannot not exceed $\epsilon^3/(300|N|^3M)$.

Define \mathcal{B} to be $\mathcal{H}_\omega \setminus (\bigcup_{i=1}^4 \mathcal{D}_i)$. If the game is not valued let the function $r^n: \mathbf{H}_\omega \rightarrow \mathbf{R}$ for player n be v_σ^n and if the game is valued then let it be the value function, already labeled r^n . Let $\hat{\sigma}$ be the candidate strategy tuple for the ϵ perfection property. We complete the proof with the more general v_σ^n , as the proof with the valued property is the same.

Let h be any history in \mathcal{B} terminating at some $s \in S$, and let a^n be an action chosen by player n with positive probability according to $\hat{\sigma}_s^n(h)$. Due to non-membership in \mathcal{D}_3 , when player n uses any action a^n not in $B^n(h)$ against the distributions $\hat{\sigma}_s^k(h)$ for $k \neq n$ the expected value of either χ^n or v_σ^n on the next stage does not differ by more than $\epsilon/5$ from what it would be with the distributions $\sigma_s^k(h)$. By non-membership in \mathcal{D}_4 this is enough for $v_\sigma^n(h) \geq j_{\hat{\sigma}}^n(h) - \epsilon$. Because a^n is not in $A_\sigma^n(h)$, the expected value of v_σ^n on the next stage resulting from using a^n is within $\epsilon/10$ of $v_\sigma^n(h)$, and due to non-membership in \mathcal{D}_2 $v_\sigma^n(h)$ does not differ from $v_{\hat{\sigma}}^n(h)$ by more than $\epsilon/5$.

Left is to show that the probability of leaving \mathcal{B} according to $\mu_{\hat{\sigma}}$ does not exceed ϵ . Because the only difference between σ and $\hat{\sigma}$ results from the use of actions in $B^n(h)$ for some $h \in \mathcal{H}_\omega$ and $n \in N$ and the probability of such an action ever being used according to μ_σ does not exceed $\epsilon^2/(35M|N|^2)$, the conclusion follows from the unlikelihood according to μ_σ of leaving \mathcal{B} . ■

QUESTION 1: *Does there exist a normal stochastic game that is perfect but doesn't have approximate equilibria?*

There are two problems with the converse direction of Theorem 1. First, although the probability of leaving the set \mathcal{B} of histories is very small if the players stick to the suggested strategies, a player could steer intentionally the play away from \mathcal{B} with unknown consequences. Second, for an ϵ equilibrium we would like to punish a player who obtains a cumulative advantage of more than ϵ . For some fixed positive ϵ how do we know that there is a positive δ so small that if a player could obtain at most a payoff advantage of δ on each stage of play then through honest selection of actions this could not accumulate to an advantage of more than ϵ over time? Under what conditions the converse of Theorem 1 is possible is an interesting topic, partly because Vieille's proof (for two-person games) uses a special version of this converse direction. In this version, however, the relation between the ϵ and δ mentioned above is dependent on the number of different situations calling for different behavior.

Define an ϵ perfect strategy tuple σ to be ϵ **self-perfect** if the function $r^n: \mathcal{H}_\omega \rightarrow \mathbf{R}$ defining the perfection property is equal to v_σ^n . A stochastic game is **self-perfect** if there exists an ϵ self-perfect strategy tuple for every positive ϵ .

QUESTION 2: *Does there exist a normal stochastic game that is perfect but not self-perfect?*

3. An example

The following is an example of a recursive two-player normal stochastic game with limit average payoffs, three non-absorbing states and three actions for each player at each state, that is not value-perfect, meaning also that the game is not valued (by Theorem 1).

Example 1: The three non-absorbing states are the set $\{s, t, u\}$. The game is symmetric with respect to the two players, Player One and Player Two. Like a tennis match, there is a state t that means “advantage” to Player One and another state u that means “advantage” to Player Two.

At the state s both players have three actions, a for “advance”, w for “wait”, and c for “check”. For Player $i = 1, 2$ the actions a , w , and c will be called a_i , w_i , and c_i . The transitions and payoffs are as follows.

If Player One chooses a_1 at state s and

- if Player Two chooses a_2 , then the play returns to the state s ,
- if Player Two chooses w_2 , then the play advances to the state t ,
- if Player Two chooses c_2 , then the game ends with certainty with a payoff of 10 to Player One and 14 to Player Two.

If Player One chooses w_1 at state s and

- if Player Two chooses a_2 , then the play advances to the state u ,
- if Player Two chooses w_2 , then the play returns to the state s , and
- if Player Two chooses c_2 , then the game ends with a payoff of 25 to Player One and $10^{-3} = 1/1000$ to Player Two.

If Player One chooses c_1 at state s and

- if Player Two chooses a_2 , then the game ends with a payoff of 14 to Player One and 10 to Player Two,
- if Player Two chooses w_2 , then the game ends with a payoff of 10^{-3} to Player One and a payoff of 25 to Player Two, and

– if Player Two chooses c_2 , then the game ends with a payoff of 10^{-3} to both players.

At the state t Player One has three actions, e_1 for “end”, r_1 for “retreat”, and f_1 for “flip”. Player Two also has three actions, n_2 for “normal”, b_2 for “first bluff” and l_2 for “second bluff”.

If Player One chooses e_1 at the state t , then no matter what Player Two does the game ends and Player One receives a payoff of 20, and

- if Player Two chooses n_2 or l_2 , then Player Two receives a payoff of 21 and
- if Player Two chooses b_2 , then Player Two receives a payoff of $21 + 10^{-3}$.

If Player One chooses r_1 at the state t and

- if Player Two chooses n_2 , then with $1/2$ probability the game ends with a payoff of 25 to Player One and 1 to Player Two and with $1/2$ probability the play returns to state s ,
- if Player Two chooses b_2 , then the game ends with a payoff of $20 + 10^{-3}$ to Player One and 10^{-3} to Player Two, and
- if Player Two chooses l_2 , then the game ends with a payoff of 20 to Player One and $7.5 + 10^{-3}$ to Player Two.

If Player One chooses f_1 at the state t and

- if Player Two chooses n_2 , then with $2/7$ probability the play moves to the state u and with $5/7$ probability the game ends with a payoff of 19.61 for Player One and 10^{-3} for Player Two,
- if Player Two chooses b_2 , then the game ends with a payoff of 20 for Player One and $\frac{40}{7} + 10^{-3}$ for Player Two, and
- if Player Two chooses l_2 , then the game ends with a payoff of $20 + 10^{-3}$ for Player One and 10^{-3} for Player Two.

At state u the situation is symmetric to that of state t , but with the roles of the players switched. Player Two has the actions e_2 , r_2 and f_2 and Player One the actions n_1 , b_1 and l_1 . For example, if n_1 and f_2 are chosen then there is a $2/7$ probability of moving to the state t and a $5/7$ probability of the game ending with a payoff of $1/1000$ for Player One and 19.61 for Player Two.

The game starts at the state s . One should interpret the “end” of the game with corresponding payoffs to be a transition to an absorbing state. (Strictly speaking, our example has 23 absorbing states, 9 following directly after either t or u , and 5 following after s . By choosing vectors in \mathbf{R}^2 that contain all the payoffs in their convex hull we could reduce the number of absorbing states to

three.) If the game never reaches an absorbing state then both players receive a payoff of 0.

Most critical to this example is the approximate value of 15 for a player at the state s . If Player One has an expected a payoff of 15 at s , then at the state t if Player Two chooses the action n_2 Player One will be indifferent between ending the game immediately with the action e_1 or moving back to s with $1/2$ probability with the action r_1 .

When analyzing this game the term **frequency** refers to the probability that an action is chosen. This is done to avoid confusion with other expressions of probability.

LEMMA 1: *From a start at state s each player can guarantee 10.4, meaning that $\chi^1(s)$ and $\chi^2(s)$ are at least 10.4.*

Proof: By symmetry, it suffices to show that Player One can guarantee 10.4. Let Player One choose the action e_1 at state t and the action n_1 at state u . At state s let Player One choose the actions a_1 , w_1 , and c_1 with the frequencies .39, .26, and .35, respectively. We need to check that with the choice of any actions for Player Two that Player One receives at least 10.4 conditioned on not returning to the state s . If Player Two chooses the action c_2 then the payoff for Player One would be at least $.39 \cdot 10 + .26 \cdot 25 = 10.4$. If Player Two chooses the action w_2 then the payoff for Player One would be at least $\frac{.39 \cdot 20}{.74}$, which is more than 10.4. If Player Two chooses the actions a_2 and r_2 then Player One would get at least $\frac{.13 + .35 \cdot 14}{.48}$, which is more than 10.4. If Player Two chooses the actions a_2 and f_2 then Player One would get at least $\frac{(2/7) \cdot .26 \cdot 20 + .35 \cdot 14}{.61}$, which is more than 10.4. If Player Two chooses the actions a_2 and e_2 then Player One would get more than 16. ■

It is now easy to confirm that $\chi^1(t) \geq 20$ and $\chi^2(u) \geq 20$ and that $\chi^1(u) \geq 5.7$ and $\chi^2(t) \geq 5.7$.

THEOREM 2: *There is no 10^{-19} value-perfect strategy pair for the game of Example 1 (implying by Theorem 1 that the game is not valued).*

Before proving Theorem 2 we need some more lemmatta, based on the assumption that there exists a 10^{-19} value-perfect strategy pair σ , with \mathcal{B} the corresponding subsets of histories and $v^i: \{s, t, u\} \rightarrow \mathbf{R}$ the value functions for the players $i = 1, 2$. (For the absorbing states the value functions are already determined.) If a player i prefers some action over another or can obtain some quantity with an action, we are referring either to the expected value of v^i on

the next stage if the preferred action was given a positive frequency by the strategies or otherwise the expected value of χ^i on the next stage.

LEMMA 2: *The probability of never reaching an absorbing state from the start of the game cannot exceed 10^{-17} .*

Proof: If this probability did exceed 10^{-17} , then there must be a history h in \mathcal{B} such that the probability of not reaching an absorbing state in the future is at least .9. Because both players can obtain a payoff of at least 5.7 from any history terminating at a non-absorbing state, the “jump” functions j_σ^i are at least $5.7 - 10^{-18}$ for all histories terminating at non-absorbing states, meaning also that v^i is at least $5.7 - 2 \cdot 10^{-18}$ for all of these states and $i = 1, 2$. Since the functions v^i represent approximately what the players receive in the future from histories in \mathcal{B} , a member of \mathcal{B} with such a low probability of future absorption would not be possible.

LEMMA 3: *From the start of the game the subset $\{t, u\}$ of states is reached with a probability of more than $2 \cdot 10^{-3}$.*

Proof:

CASE 1; $v^i(s) \geq 12.6$ FOR BOTH PLAYERS: It follows directly from the fact that the sums over the two players from all absorbing payoffs following directly after the state s (without first reaching either t or u) never exceed 25.001 and this sum for all absorbing payoffs of the game never exceeds $41 + 10^{-3}$.

CASE 2; $v^i(s) < 12.6$ FOR SOME PLAYER i : By symmetry we assume that $v^1(s) < 12.6$. The frequency given to a_2 at any history in \mathcal{B} terminating at s does not exceed $11/12$, otherwise Player One could receive at least 12.66 on this stage by playing c_1 .

Let h_0 be the initial history at the first stage of play at the state s , which necessarily belongs to \mathcal{B} .

CASE 2a; PLAYER ONE CHOOSES a_1 WITH POSITIVE FREQUENCY AT h_0 : From above we know that $\sigma_s^2(h_0)(c_2) + \sigma_s^2(h_0)(w_2) \geq 1/12$. If $\sigma_s^2(h_0)(c_2)$ were not at least twice that of $\sigma_s^2(h_0)(w_2)$ then $w_\sigma^{v^1}(h_0)(a_1)$ would be at least 10^{-18} more than $v^1(s) < 12.6$, a contradiction. Likewise $25\sigma_s^2(h_0)(w_2) \geq \sigma_s^2(h_0)(c_2)$, since otherwise $w_\sigma^{v^1}(h_0)(a_1) \geq v^1(s) - 10^{-18}$ would imply that $v^1(s)$ is less than $10 + 10^{-15}$, a contradiction to Lemma 1. With Player Two choosing the action c_2 with positive frequency it is necessary that $\sigma_s^1(h_0)(a_1) \geq 7/10$. But we have shown above that $\sigma_s^2(h_0)(w_2) \geq 1/312$. This implies that t is reached on the second stage with a probability of at least $\frac{1}{312} \frac{7}{10} > 1/500$.

CASE 2b; PLAYER ONE DOES NOT CHOOSE a_1 WITH POSITIVE FREQUENCY AT h_0 : Clearly Player One chooses w_1 with positive frequency at h_0 , since otherwise Player One would choose c_1 with certainty and Player Two would choose w_2 with certainty. We know from Lemma 1 that c_2 is not chosen with positive frequency. With a payoff of at least $10.4 - 10^{-18}$ for Player One from the combination of a_1 with a_2 (a return to state s), to prevent Player One from getting at least $12.6 + 10^{-3}$ we must assume that Player Two chooses the action a_2 with a frequency of at least .6 at the history h_0 . If we assume that Player One chooses the action c_1 with a frequency of at least .98 at the history h_0 , then Player Two would have chosen the action w_2 with certainty. Therefore we must assume that the probability of reaching u on the second stage must be at least $.6/50 = .012$. ■

LEMMA 4: *If there are members of \mathcal{B} terminating at both t and u then at any such histories the corresponding action r_i is chosen with positive frequency.*

Proof: For the sake of contradiction we assume that the action r_1 is not chosen with positive frequency at some history h in \mathcal{B} terminating at t . By comparing the actions l_2 and b_2 we can assume also that l_2 is also not chosen with positive frequency at h . If Player Two preferred the action b_2 by more than 10^{-18} over the action n_2 , then only b_2 would have been chosen and then indeed Player One would prefer the action r_1 over the others by more than 10^{-18} . Furthermore, if Player Two placed all but 10^{-16} frequency on the action b_2 then the same would hold. To prevent such a preference for the action b_2 over the action n_2 it would be necessary for $v^2(u)$ to be at least $20 + 10^{-4}$ and that Player One chooses the action f_1 with positive frequency.

CASE 1; $v^1(u) \leq 20.9$: To prevent Player One preferring the action e_1 over the action f_1 by more than 10^{-18} at the history h , it would be necessary that Player Two chooses the action l_2 with positive frequency at h . But as argued above, without a positive frequency for the action r_1 this is not possible.

CASE 2; $v^1(u) \geq 20.9$: It would be necessary that Player Two chooses the action e_2 with positive frequency for all histories in \mathcal{B} terminating at u . But this implies that $v^2(u)$ is no more than $20 + 10^{-18}$, a contradiction to the above assumption. ■

LEMMA 5: *If Player Two chooses b_2 with positive frequency at a history in \mathcal{B} terminating at t it is necessary that Player One also chooses r_1 at this same*

history with a frequency less than $1/5000$. If additionally Player Two chooses both b_2 and l_2 with positive frequency then at this same history Player One chooses the action f_1 with a frequency less than $1/10000$. The corresponding symmetric statement from switching the players and the states holds.

Proof: From only the actions e_1 and f_1 the advantage for Player Two by choosing b_2 over the action n_2 does not exceed $1/1000$. Since Player Two gets at least 5.7 from the combination of r_1 with n_2 , if Player One chose r_1 with a probability of at least $1/5000$, Player Two would prefer by more than 10^{-18} the action n_2 over the action b_2 . Additionally, if Player One chose the action f_1 with a frequency exceeding $1/10000$, Player Two would prefer by more than 10^{-18} the action b_2 over the action l_2 . (Player Two would lose no more than $\frac{15}{10000}$ from the combination of b_2 with r_1 , but then would gain at least $\frac{4}{7000} + \frac{.9997}{1000}$ from the combinations of b_2 with e_1 and f_1 .) ■

LEMMA 6: If $v^1(s) \leq 14.9$, $v^1(u) \leq 20.9$ and there is some history in \mathcal{B} terminating at t then $v^2(t)$ is at least 20.995. The corresponding symmetric statement from switching the players and the states holds.

Proof: For every history in \mathcal{B} terminating at t the above assumptions imply that the combination of r_1 with n_2 gives no more than 19.95 to Player One and the combination of f_1 with n_2 gives no more than 19.98 to Player One (with respect to the expected value of v^1 on the next stage).

First, at such a history one can assume that b_2 was chosen with positive frequency. Assume the contrary. If Player Two chose n_2 with a frequency of at least 10^{-16} then Player One would prefer e_1 over r_1 by more than 10^{-18} . Without Player One choosing r_1 with positive frequency and without Player Two choosing b_2 with positive frequency, the only way to prevent Player Two from only choosing the action n_2 (which would lead quickly to a contradiction) would be that Player One chooses e_1 with a frequency of at least $1 - 10^{-18}$. But this would imply the conclusion of the lemma, and therefore we can assume that the frequency for l_2 was at least $1 - 10^{-16}$. But this leads directly to a contradiction (as Player One would respond by choosing only f_1). Therefore we can assume that b_2 was chosen indeed with positive frequency.

Second, at such a history one can assume that Player Two chose l_2 with positive frequency. Suppose the contrary. We can assume that Player One didn't choose f_1 with positive frequency, since the only way to prevent a 10^{-18} preference for e_1 over f_1 would be if all but 10^{-14} frequency went to the action

b_2 (also leading to a contradiction). But with no weight given to the action f_1 the only way to prevent Player Two from choosing only the action n_2 (which would lead to a contradiction) would be for Player One to choose e_1 with a frequency of at least $1 - 1/5600$. The result of such a behavior would also imply the conclusion of the lemma.

With both b_2 and l_2 chosen with positive frequency, Lemma 5 completes the proof (for example by looking at Player Two's option to choose b_2). ■

LEMMA 7: *From any member of \mathcal{B} terminating at t the action n_2 is chosen with a frequency of at least $5 \cdot 10^{-5}$. If additionally $v^2(s) \leq 14 + 10^{-5}$ then the action f_1 is chosen with a frequency of at least $4 \cdot 10^{-9}$. If Player One does not choose r_1 with positive frequency at a member of β terminating at t then the action f_1 is chosen at this history with a frequency of at least $5 \cdot 10^{-5}$. The corresponding symmetric statement from switching the players and the states holds.*

Proof: First we show that Player Two chooses the action n_2 with a frequency of at least $5 \cdot 10^{-5}$. Suppose the contrary. The first consequence is that Player One does not choose e_1 with positive frequency, since either the action r_1 or the action f_1 would be preferable by a quantity of at least 10^{-4} . Player Two must give positive frequency to both b_2 and l_2 , in the first case to prevent Player One from choosing only the action f_1 and in the second case to prevent Player One from choosing only the action r_1 (as both would result in a contradiction). But then by Lemma 5 Player One must choose e_1 with positive frequency (which would be a contradiction).

A frequency for e_1 above $1 - 4 \cdot 10^{-5}$ would imply that Player Two chooses only b_2 , which would lead to a contradiction, so the frequencies for r_1 and f_1 add up to at least $4 \cdot 10^{-5}$. Given $v^2(s) \leq 14 + 10^{-5}$, Player Two prefers the combination of l_2 with r_1 by at least $8 \cdot 10^{-4}$ over the combination of n_2 with r_1 . To get Player Two to choose n_2 with positive probability (which must hold by the above) it would be necessary for Player One to choose f_1 with a frequency of at least $4 \cdot 10^{-9}$.

Lastly, if Player One did not choose r_1 with positive frequency then one must conclude that Player Two did not choose l_2 with positive frequency (since b_2 would be a much preferable action). Since Player Two can get no more than $51/7$ from the combination of n_2 and f_1 , it is necessary that Player One chooses f_1 with a frequency of at least $5 \cdot 10^{-4}$ to prevent Player Two from preferring b_2 over n_2 by more than 10^{-18} . ■

Proof of Theorem 2: We separate the proof into three cases, determined by the use of the actions c_i .

CASE I; NEITHER PLAYER i CHOOSES THE ACTION c_i AT ANY HISTORY IN \mathcal{B} TERMINATING AT s :

CASE IA; FOR ONLY ONE OF THE STATES t OR u IS THERE A HISTORY IN \mathcal{B} TERMINATING AT THIS STATE:

By symmetry we can assume that state t is reached directly from s (by a combination of the actions a_1 with w_2) with a probability of at least $1 - 10^{-17}$ and that $v^1(s) \geq 20 - 10^{-15}$. Due to Lemma 7 (implying minimal frequencies for the choices of f_1 and n_2 at the state t) and the assumption on not reaching u with histories in \mathcal{B} , we can assume that the probability of reaching a history terminating at t where Player One does not choose r_1 with positive probability does not exceed 10^{-8} . At any history in \mathcal{B} terminating at t where Player One chooses the action r_1 with positive frequency, the combination of r_1 with n_2 gives at least 22 for Player One and due to Lemma 7 (implying a minimal frequency for n_2) the action r_1 dominates the action e_1 by more than 10^{-18} . With Player One not choosing e_1 with positive frequency at such a history, no matter what Player Two does, we must assume that $v^2(s)$ is no more than 8, a contradiction to Lemma 1.

CASE IB; THERE ARE HISTORIES IN \mathcal{B} TERMINATING AT BOTH t AND u : By symmetry we can assume that t is reached directly from s (by a combination of a_1 and w_2) with a probability of at least .499.

CASE IBi; $v^1(s) \geq 15.1$: According to Lemma 4 we know that Player One chooses r_1 with positive frequency at any history in \mathcal{B} terminating at t . By Lemma 7 (implying a minimal frequency for n_2) and by comparing the action r_1 with e_1 , we can conclude that Player One is not choosing e_1 with positive frequency at such a history, and therefore the payoff for Player Two conditioned on reaching t and not returning to s cannot exceed 7.6.

Now we consider any history in \mathcal{B} terminating at s such that both a_1 and w_2 are chosen with positive frequency. Since Player One does not choose the action c_1 with positive frequency nor e_1 at any history in \mathcal{B} terminating at t , if Player One chooses a_1 with a frequency of at least $3 \cdot 10^{-19}$ then by $v^2(t) < v^2(s) - 1.2$ (from Lemma 1 and the conclusion of the last paragraph) Player Two would not choose the action w_2 with positive frequency. So with Player One choosing a_1 with a frequency of no more than $3 \cdot 10^{-19}$, Player Two can receive a payoff of at least $20 - 3 \cdot 10^{-18}$ by choosing the action a_2 instead, and therefore $v^2(s) \geq$

$20 - 3 \cdot 10^{-18}$. Since t is reached with a probability of at least .499 and Player One does not choose either e_1 or c_1 with positive probability (at any histories in \mathcal{B}), the amount Player Two gets from the state s , namely $v^2(s)$, could not be more than 19, a contradiction.

CASE IBii; $v^1(s) \leq 15.1$: Since Player One receives at least $20 - 10^{-18}$ from the state t and this state is reached from s with a probability of at least .499, $v^1(u)$ can be no more than 10.25, at least .15 smaller than the guaranteed payoff for Player One at the state s . Furthermore, there must be a history in \mathcal{B} terminating at s where Player One chooses w_1 with positive frequency, since otherwise $v^1(t) \geq 20 - 10^{-18}$ would make $v^1(s) \leq 15.1$ impossible. Whenever Player One chooses w_1 with positive frequency at a history in \mathcal{B} terminating at s , it would be necessary that Player Two chooses w_2 with a frequency of at least $1 - 10^{-17}$ (since otherwise by $v^1(u) \leq 10.25$ the choice of w_1 would result in an expected payoff for Player One less than $v^1(u) - 10^{-18}$). But then Player One could choose a_1 for a guaranteed payoff of at least $20 - 10^{-15}$, a contradiction to the assumption $v^1(s) \leq 15.1$.

CASE II; PLAYER i CHOOSES THE ACTION c_i AT SOME HISTORY IN \mathcal{B} TERMINATING AT s , BUT PLAYER $j \neq i$ DOES NOT CHOOSE THE ACTION c_j AT ANY HISTORY IN \mathcal{B} TERMINATING AT s : By symmetry, we assume that it is Player One who chooses the action c_1 . Since Player One can get no more than 14 with the action c_1 , we have that $v^1(s) \leq 14 + 10^{-18}$. This implies for every history h in \mathcal{B} terminating at s that the frequency for a_2 must be at least $3/5$, otherwise Player One would prefer to choose the action a_1 , even when punished in the event that Player Two had chosen a_2 at the same time. Likewise, the frequency for a_2 must be at least $\frac{10.399}{14} > .74$ whenever Player One is choosing c_1 with positive frequency.

Next we know that there is some history in \mathcal{B} terminating at t . To contradict this claim it would be necessary by Lemma 3 that the state u is reached with a probability of at least $\frac{1}{501}$. But then Lemma 7 (its symmetric statement) would imply that t is reached from u with a probability of at least $2 \cdot 10^{-13}$, implying that indeed t is reached with a probability of at least 10^{-17} .

CASE IIA; THERE IS NO HISTORY IN \mathcal{B} TERMINATING AT u OR THERE IS NO HISTORY IN \mathcal{B} TERMINATING AT s WHERE PLAYER ONE CHOOSES w_1 WITH POSITIVE FREQUENCY: Assume first that there is no history in \mathcal{B} terminating at u . Because Player Two chooses a_2 with a frequency of at least $3/5$, there must be a subset \mathcal{B}' of \mathcal{B} such that $\mathcal{B} \setminus \mathcal{B}'$ is reached with a probability not

exceeding 10^{-9} and in every history in \mathcal{B}' Player One does not choose w_1 with a frequency exceeding 10^{-9} . If, however, Player One does not choose w_1 with positive frequency at any history in \mathcal{B} , proceed with the assumption that \mathcal{B}' is equal to \mathcal{B} .

Consider any history h in \mathcal{B}' where Player One chooses c_1 with positive frequency and Player Two chooses w_2 with positive frequency. Because Player One chooses w_1 with at most a frequency of 10^{-9} and $v^2(s) \geq 10.4$, it follows that Player One is choosing a_1 with a frequency of at least $1 - 2 \cdot 10^{-8}$ (since otherwise Player Two would lose too much from the action a_2 and choose only the action w_2 , implying of course that Player One could not have opted for the action c_1). As Player One receives at least $v^1(s) + 5\sigma_s^2(h)(w_2)$ from the action a_1 (as $v^1(t) \geq 20 - 10^{-18}$), the frequency $\sigma_s^2(h)(w_2)$ given to w_2 by Player Two at h cannot exceed 10^{-19} . Since this is true for all h in \mathcal{B}' where c_1 and w_2 are chosen, as a consequence the combination of c_1 with w_2 plays no significant role in the expected payoffs of the game and the two-dimensional vector $(v^1(s), v^2(s))$ is within a Euclidean distance of 10^{-5} from the convex combination of $(14, 10)$ and $(v^1(t), v^2(t))$. But this is not possible, since $v^1(t) \geq 20 - 10^{-18}$ and $v^1(s) \leq 14 + 10^{-18}$, which would imply that $v^2(s)$ is less than 10.1, a contradiction to Lemma 1.

CASE IIB; THERE IS SOME HISTORY IN \mathcal{B} TERMINATING AT u AND THERE IS SOME HISTORY IN \mathcal{B} WHERE PLAYER ONE CHOOSES w_1 WITH POSITIVE FREQUENCY: Since Player Two chooses the action a_2 with a frequency of at least $3/5$ at any history in \mathcal{B} terminating at s , we must assume that when Player One does choose w_1 with positive frequency at some history in \mathcal{B} terminating at s the quantities $v^1(u)$ and $v^1(s)$ must be within 10^{-18} of each other, implying that $v^1(u) \leq 14 + 10^{-17}$. By Lemma 6 we have $v^2(t) \geq 20.995$. By Lemma 7 at all histories in \mathcal{B} terminating at u , Player Two chooses f_2 with a positive frequency of at least $4 \cdot 10^{-9}$. With $\frac{5}{7} \cdot 19.61 + \frac{2}{7} \cdot 20.995 = 20 + \frac{4}{700}$ there would be an advantage of at least $2 \cdot 10^{-11}$ to Player Two by playing f_2 instead of e_2 , and therefore e_2 is never chosen at any history in \mathcal{B} terminating at u , implying that $v^1(u) \leq 7.9$. But by the above this implies that $v^1(s)$ is also below 8, a contradiction to Lemma 1.

CASE III; FOR BOTH $i = 1, 2$ PLAYER i CHOOSES THE ACTION c_i AT SOME HISTORY IN \mathcal{B} TERMINATING AT s : It is necessary that $v^i(s) \leq 14 + 10^{-18}$ for both $i = 1, 2$, since otherwise neither c_1 nor c_2 would have been used. Lemmata 3 and 7 imply that there are histories in \mathcal{B} terminating at t and at u . Also from $v^i(s) \leq 14 + 10^{-18}$ for both $i = 1, 2$ we conclude that both $v^1(t)$ and $v^2(u)$ do

not exceed 20.1.

Furthermore, Player One must choose a_1 with a frequency of at least $1/50$ at any history in \mathcal{B} terminating at s , and the same must hold for Player Two and a_2 . On the one hand, if Player One chose c_1 with a frequency of at least .28 and also a_1 with less than $1/50$ then Player Two could get at least 14.1 from responding with the action w_2 (even if he is punished for choosing this action). On the other hand, if Player One chose w_1 with a frequency of at least .7 then Player Two could get at least 14.1 by choosing the action a_2 .

Also by Lemma 3 there is some history h in \mathcal{B} terminating at s where the combination w_i and a_k are played, with $i \neq k$. Without loss of generality we can assume that the actions a_1 and w_2 were played together at some history h in \mathcal{B} terminating at s .

CASE IIIA; THE ACTION w_1 WAS NOT CHOSEN WITH POSITIVE PROBABILITY AT THE HISTORY h : With a frequency of at least $1/50$ given to a_1 , if $v^2(t)$ were at least $14 + 10^{-16}$ then $w^{v^2}(h)(w_2)$ would be at least $14 + 2 \cdot 10^{-18}$, which would contradict the main assumptions of Case III. From Lemma 6 we must conclude that $v^1(u)$ is at least 20.995. Therefore from Lemma 7 (establishing minimal frequencies for the actions n_2 and f_1) Player One chooses the action f_1 with positive frequency and prefers it over the action e_1 by more than 10^{-18} at any history in \mathcal{B} terminating at t , implying that $v^2(t) < v^2(s) - 1$. Due to the minimal frequency of $1/50$ given to a_1 and Player Two's positive frequency for w_2 , we must presume at the history h that Player One also chose c_1 with a frequency of at least $1/1000$ (as otherwise the choice of w_2 would result in a payoff of no more than $v^2(s) - 1/200$). Since Player One is not choosing w_1 with positive frequency at the history h , Player Two could not have chosen the action a_2 at h with positive frequency, as then $w_\sigma^{v^2}(h)(a_2) \leq \frac{10}{1000} + \frac{999}{1000}v^2(s) \leq v^2(s) - 10^{-4}$ (from $v^2(s) \geq 10.4 - 10^{-18}$). But this contradicts our argument that Player One does choose c_1 with positive frequency at the history h .

CASE IIIB; THE ACTION w_1 WAS CHOSEN WITH POSITIVE PROBABILITY AT THE HISTORY h : Since both a_1 and a_2 were chosen with a frequency of at least $1/50$, we must assume that $v^1(u)$ and $v^2(t)$ are both no more than $14 + 10^{-16}$ (since otherwise the actions a_1 and a_2 would not have been chosen, in preference for w_1 and w_2). On the other hand, by Lemma 6 we must assume that both $v^1(u)$ and $v^2(t)$ are at least 20.995, a contradiction. ■

We know of at least two ϵ equilibrium strategies to this game.

For the first approximate equilibrium, at all visits to s let the players alternate between playing a_1 with w_2 and playing a_2 with w_1 . When at the state t

let Player Two choose only the action n_2 and let Player One choose e_1 with a frequency of $9/21$ and r_1 with a frequency of $12/21$. Let the players act symmetrically at the state u . The future expected payoffs for their visits to the state s will alternate between $(15, 20)$ and $(20, 15)$, which will imply that Player One will be indifferent between her two actions at the state t and the same holds for Player Two at the state u . It is easy to check that neither player has any motivation to deviate.

The second equilibrium corresponds to the Vieille proof, and it is not so easy to find. Almost all of the time both players perform together the actions a_1 and a_2 at the state s . At the state t , Player Two will choose n_2 with certainty and Player One will choose the actions r_1 and e_1 with frequencies $7/13$ and $6/13$, respectively. At the state u , Player One will choose the action n_1 with certainty and Player Two will choose the actions r_2 and e_2 with frequencies $6/13$ and $7/13$, respectively. Every time the play is at the state s the players count the number of times the play has been consecutively in s . Let N be any natural number greater than $100/\epsilon$. If it is the N th consecutive visit to s or something less, then Player Two plays a_2 with certainty. In those initial N stages Player One chooses at each stage the action w_1 with a positive frequency p which satisfies $(1 - p)^N = 5/6$ (and otherwise a_1 is chosen). If the players reach the $N + 1$ st consecutive visit to s , then Player One will play a_1 with certainty and Player Two will choose w_2 with the frequency of a very small quantity δ (for example $\delta < \epsilon/40$ suffices) and c_2 also with the same frequency of δ . They continue in this way until either w_2 or c_2 is chosen. One can calculate that the expected payoffs for both players at a first visit to s will be 15. Player One gets always an expected payoff of 15 at states s and u , and Player Two receives an expected payoff of 14 from the combination of w_2 and a_1 . The expected payoff for Player Two at the state s varies within each extended visit to that state, starting at 15 and falling to 14 if Player One fails to choose w_1 . With δ sufficiently small, Player One cannot gain more than ϵ by choosing w_1 when Player Two is choosing w_2 or c_2 with the very small frequency δ .

4. Conclusion: Countably many states

The discovery of the above example originated from our curiosity concerning the existence of approximate equilibria for two-player normal stochastic games with countably many states. Our pessimism concerning the existence of approximate equilibria for these games has its origin in the complexity of the ϵ equilibrium strategies of Vieille's proof. If a game has finitely many states and there is some

strategy pair such that the expected number of visits to any non-absorbing state is finite, then almost surely an absorbing state will be reached (something that may not be true if there are infinitely many non-absorbing states). This was used critically in Vieille's proof to show that payoffs (averaged over the stages) were converging almost surely. If there is a proof of the existence of approximate equilibria for two-person games with countably many states, one could expect that its application to games with finitely many states would deliver ϵ equilibria with much faster rates of payoff convergence than that obtained from the Vieille proof.

There is an additional equilibrium concept more problematic when there are infinitely many states, that of uniformity. An ϵ equilibrium for a limit average stochastic game is **uniform** if for some sufficiently large N for all $n \geq N$ it defines ϵ equilibria for the truncated games that end at the n th stage (and where the payoffs are determined by the averages over the n stages). Vieille's proof is that of the existence of uniform ϵ equilibria for every positive ϵ . We choose to consider a class of games for which the existence of approximate equilibria implies that there are uniform ϵ equilibria for every $\epsilon > 0$. Assume a normal stochastic game is recursive with fixed payoffs for all players at all absorbing states and assume that some player can obtain a positive $w > 0$ from a start at any non-absorbing state. With any ϵ equilibrium the probability that an absorbing state is not reached from the start of the game cannot exceed ϵ/w , and therefore there must be a stage K such that the probability of absorption before reaching the stage K is greater than $1 - 2\epsilon/w$. Assuming that positive ϵ is less than 1 and letting N be greater than K/ϵ , the original ϵ equilibrium defines $\epsilon M + 2\epsilon M/w$ equilibria for all games of n stages with $n \geq N$ (where M was the bound on the maximal difference between any two payoffs).

A strategy for finding a counter-example for countably many states using the above class of recursive games could be the following. Construct an infinite sequence of recursive games $\Gamma_0, \Gamma_1, \dots$ with increasing finite sets $S_0 \subseteq S_1 \subseteq \dots$ of non-absorbing states such that for every $i \geq 0$ and $j \geq i$ the actions in S_i are the same for all games Γ_j , all payoffs for both players from absorbing states are greater than one, and if a is an action tuple at $s \in S_i$ then $p(i)_a^s(t) = p(j)_a^s(t)$ for all $j \geq i$ and $t \in S_i$, (with $p(k)$ the transition laws for the game Γ_k). Construct a game played on a countable state space ($\bigcup_{i=0}^{\infty} S_i$ unioned with the absorbing states from all the games Γ_i) by having the game start at $\hat{s} \in S_0$, define the non-absorbing states on the i th stage to be the set S_i , and declare that absorption occurs if an absorbing state of the game Γ_i has been reached

on stage i . Furthermore, give both players the ability to force the game to the set of absorbing states with certainty from any start at a non-absorbing state of the new countable state space. Desirable may be games Γ_i such that with large i the approximate equilibrium behavior of Γ_i keeps the non-absorbing play most of the time close to the set S_0 and the minimal number of stages necessary to reach an absorbing state in the game Γ_i starting from any $s_0 \in S_0$ goes to infinity as i goes to infinity.

References

- [1] D. Blackwell and T. S. Ferguson, *The Big Match*, *Annals of Mathematical Science* **39** (1968), 159–163.
- [2] J. Flesch, F. Thuijsman and O. J. Vrieze, *Cyclic Markov equilibria in a cubic game*, *International Journal of Game Theory* **26** (1997), 303–314.
- [3] A. Maitra and W. Sudderth, *An operator solution of stochastic games*, *Israel Journal of Mathematics* **78** (1991), 33–49.
- [4] D. Martin, *The determinacy of Blackwell games*, *The Journal of Symbolic Logic* **63** (1998), 1565–1581.
- [5] J. F. Mertens and A. Neyman, *Stochastic games*, *International Journal of Game Theory* **10** (1981), 53–66.
- [6] J. Nash, *Equilibrium points in n -person games*, *Proceedings of the National Academy of Sciences of the United States of America* **36** (1950), 48–49.
- [7] L. S. Shapley, *Stochastic games*, *Proceedings of the National Academy of Sciences of the United States of America* **39** (1953), 1095–1100.
- [8] N. Vieille, *Two-player stochastic games I: A reduction*, *Israel Journal of Mathematics* **119** (2000), 55–91; *Two-player stochastic games II: The case of recursive games*, *Israel Journal of Mathematics* **119** (2000), 92–126; *Small perturbations and stochastic games*, *Israel Journal of Mathematics* **119** (2000), 127–142.