

Discussion Paper No. 1191

**Nash Equilibrium and the Evolution
of Preferences**

by

Jeffrey C. Ely
and
Okan Yilankaya

Northwestern University

August 1997

Nash Equilibrium and the Evolution of Preferences *

Jeffrey C. Ely[†] Okan Yilankaya[‡]

August 4, 1997

Abstract

A population of players of players is randomly matched to play a normal form game G . The payoffs in this game represent the fitness associated with the various outcomes. Each individual has preferences over the outcomes in the game and chooses an optimal action with respect to those preferences. However, these preferences needn't coincide with the fitness payoffs. When evolution selects individuals on the basis of the fitness of the actions they choose, the distribution of aggregate play must be a Nash equilibrium of G . Weak additional assumptions on the evolutionary process imply perfect equilibrium.

*Department of Economics, Northwestern University. 2003 Sheridan Road, Evanston IL 60208. We thank (but exonerate) Eddie Dekel, who contributed substantially to the ideas in this paper.

[†]ely@nwu.edu

[‡]okan-y@nwu.edu

1 Introduction

Economic models built around rational self-interested agents are rarely, if ever, accurate as literal descriptions of the environments they intend to capture. Agents' objectives may differ from those attributed to them, and even when they coincide agents may not have the sophistication necessary to choose actions which best achieve those objectives. The "as if" viewpoint is a defense of economic theory based on the following argument. Typically, what is at stake in the economic environments is important for survival as a player in that environment (for example, profits in a market context). Therefore, regardless of the actual *motives* of real-world agents, an essentially Darwinian mechanism should eventually imply that their *behavior* is consistent with optimization (else they would not have survived.) It should appear to an outside observer who is agnostic about the true decision-making process "as if" it were the outcome of strategically sophisticated interaction among optimizing agents.

Game theorists have attempted to formalize one aspect of the viewpoint using models of evolutionary equilibrium. The agents in these models are not rational utility-maximizers, but rather are genetically programmed to play particular actions. These agents interact with one another over time, and evolution selects in favor of those agents whose pre-programmed actions happen to be optimal against the (distribution of) actions of other agents. Evolutionary equilibrium, a situation in which every surviving agent uses an optimal action, provides some support for as if: the distribution of actions in an evolutionary equilibrium must be a Nash equilibrium. This is true for virtually all formalizations of evolutionary equilibrium, e.g. the ESS concept (Maynard Smith 1982) and its variants, as well as dynamic versions such as the replicator dynamics (Taylor and Jonker 1978) (See Weibull (1995) for a survey).

There are drawbacks. First of all, while solution concepts based on evolutionary ideas provide predictions which are consistent with Nash equilibrium, they sometimes make no prediction at all (ESS can fail to exist for some games, the replicator dynamics can fail to converge). Secondly, (and it can be argued that these two points are closely linked) they tend to bear too much resemblance to biological models and too little to economics. Even those who are sympathetic to the idea of "bounded rationality" can be skeptical that this extreme behavioral assumption is any more convincing as a model of economic agents.

In this paper we propose an alternative approach to the as if argument. We start by specifying an n -player game G with action sets $(A_i)_{i=1}^n$, and payoff functions $(\pi_i)_{i=1}^n$. As in the standard evolutionary framework, we interpret these payoff functions as representing fitness, and we imagine a population of individuals who are repeatedly randomly matched to play G . Unlike the standard

framework, agents in our world are rational decision-makers. They have preferences over outcomes and they form conjectures about the behavior of other agents. Based on these they make choices which are optimal given their preferences. In our interpretation, however, the term “rational” implies nothing more than this. In particular, it imposes no constraint on the preferences governing these choices (however we do assume that they satisfy the standard von Neumann-Morgenstern axioms.)¹ Which preferences are represented in the population is to be determined endogenously by the evolutionary process.

Slightly more formally, each of n populations (one for each player-role in G) is characterized by a distribution μ^i over the set Θ_i of possible utility functions. An individual is drawn independently from each population according to the product probability $\mu = \prod_1^n \mu^i$ to choose actions in G . Selected individuals know their own utility functions, have beliefs about their opponents’ play and choose actions to maximize expected utility.

This interaction can be summarized as an n -player game of incomplete information $\Gamma(\mu)$ in which the prior distribution over “type”-profiles is μ , each player observes his realized type, and chooses an action in A_i .

We assume that play is described by an equilibrium of $\Gamma(\mu)$. That is, each individual has a correct belief about the distribution of actions he will face and chooses an action which is a best reply to this distribution (according to his own preferences). This equilibrium determines an aggregate distribution of action-profiles, which to an outside observer appears as the outcome of the underlying game G .

The preferences that are “fit” are those that induce choices that are successful relative to the objective payoffs $(\pi_i)_{i=1}^n$. The population distribution of preferences evolves as those that are more fit grow in representation relative to those that are less fit. An evolutionary equilibrium is then a distribution of preferences μ and an equilibrium of $\Gamma(\mu)$ such that all individuals are equally successful relative to $(\pi_i)_{i=1}^n$.

Our question is whether play in an evolutionary equilibrium will appear to an outside observer as if it were the outcome of Nash equilibrium play by agents whose preferences were actually given by $(\pi_i)_{i=1}^n$.

¹We do not mean to advance the position that rationality implies an expected utility representation of preferences. We assume such a representation because it simplifies the analysis and allows us to focus on the question at hand, namely the evolution of preferences over outcomes of *strategic interaction*. It is indeed interesting to explore the parallel question of whether evolution should imply that behavior is consistent with expected utility maximization. This is a question about the evolution of preferences over *uncertain* outcomes, hence outside the narrow scope of this paper. Papers which have examined this question include Robson (1996b), Robson (1996a) and To (1995).

To answer this question, we propose a stability criterion for preferences based on the type of evolution discussed above. In the spirit of “static” concepts of evolutionary stability (such as ESS), our criterion is intended to capture the effects of mutation and natural selection, while avoiding an explicit model of the evolutionary process. Under our definition, a set S of preference distributions is stable if there is a set U of neighboring distributions such that starting anywhere within U , evolution must result in a return to S . The neighboring distributions are interpreted as the set of possible outcomes of a process of mutation. The “paths” of the evolutionary process are modeled abstractly as *selection sequences*: sequences of distributions which satisfy a standard “fitness monotonicity” property.

A set of outcomes in G is *supported by stable preferences* if those outcomes can be obtained as the distributions of play within a stable set of preference distributions. We show that every game G has a non-empty set of outcomes that are supported by stable preferences. We take this to be an important advantage of the present approach over the standard models built around evolution of strategies, which can fail to generate solutions in many games. By assuming nothing more than monotonicity in selection sequences, we prove that outcomes which are supported by stable preferences must correspond to Nash equilibrium distributions of G . Thus, our model formalizes the argument in favor of the *as if* viewpoint. Finally, by imposing some weak additional assumptions on selection sequences, we obtain an equilibrium refinement: only trembling-hand perfect equilibria can be supported by stable preferences.

2 A Model

We start with an n -player normal form game G with finite action sets A_i , $i \in \{1, \dots, n\}$ and payoff function $\pi : A \rightarrow \mathbf{R}^n$, where $A = \prod_n A_i$. We view π as representing the “true” objective payoffs, or fitnesses. A player’s survival is dependent upon his success in the game as evaluated by π . Let Δ represent the set of probability distributions on A , i.e. the set of *outcomes* in G , and $E \subset \Delta$ those distributions arising from Nash equilibria of G .

We follow the standard approach to evolutionary equilibrium selection by supposing there are n populations of players, and a process which randomly selects an individual from each population to play G . We depart from the standard approach by assuming these individuals have preferences over outcomes in G , and choose actions optimally in response to beliefs about the play of their selected opponents. However, these preferences are not necessarily represented by π .

Let $\Theta_i = [0, 1]^{|A_i|}$ be the set of possible von Neumann-Morgenstern payoff functions on A . Notice that with this specification, each preference ordering

is represented by a continuum of distinct, but equivalent, payoff functions. For example, Θ_i contains a continuum of affine transformations of π . This equivalence class of preferences will play an important role and will be denoted $\tilde{\pi}$. The set of possible n -vectors θ of payoff functions is $\Theta = \prod_1^n \Theta_i$, and can be thought of as the set of all *games* with action set A .

The environment will be characterized by a product probability measure on Θ representing the current distributions of preferences in each of the n populations.² It simplifies some arguments to assume that these distributions are non-atomic. Let $\mathcal{P}(\Theta)$ be the set of all non-atomic probability measures μ on Θ such that $\mu = \mu^1 \times \mu^2 \times \dots \times \mu^n$. Denote by $C(\mu)$ the support of the preference distribution μ .³ The matching process selects individuals from population i according to the distribution μ^i , independently of the players drawn from other populations.

This interaction can be described as an n -player Bayesian game $\Gamma(\mu)$ in which the set of possible states of the world is Θ and the common prior distribution is μ . We are going to assume that aggregate play among the populations is an equilibrium of this game. That is, we will assume that each individual, upon being selected to play, will have correct beliefs about the distribution over his opponents' play and will choose an action that is a best-reply to this belief, given his own preferences.

While it is not a part of our formal model, we view equilibrium as arising from a process of learning which operates much faster than the evolutionary process we seek to model. To be specific, our model will describe the evolution of the type distribution μ . We suppose that whenever a new distribution ν arises as a consequence of evolutionary forces, the learning process always reaches equilibrium with respect to $\Gamma(\nu)$ before subsequent evolution proceeds.⁴

A pure strategy profile in $\Gamma(\mu)$ is a measurable function $\sigma : \Theta \rightarrow A$ specifying an action profile in G for every possible payoff profile. Let Σ be the set of

²Throughout the paper, the domain of the preference distribution is the Borel σ -algebra of subsets of Θ

³The support of a probability measure is the smallest closed set which has measure 1. Such a set always exists for Borel measures on subsets of \mathbf{R}^n

⁴We find it perfectly natural to assume that evolution proceeds much more slowly than learning, but clearly our assumption is extreme. An interesting line of development for this research would be to explicitly embed a model of learning into our evolutionary framework. It would then be possible to ask whether our assumption is justified by a model in which the relative rates of learning and evolution are somehow parameterized and the appropriate limit is taken. A common result of dynamic models of equilibrium selection is that (the implications of) extreme assumptions are not necessarily borne out by the limits of "interior" models. Papers with such lessons include Binmore and Samuelson (1997), Binmore, Samuelson, and Vaughn (1992), Ely (1996), Kandori, Mailath, and Rob (1993), Robson and Vega-Redondo (1996), Sandholm and Pauzner (1997) and Young (1993).

pure strategy profiles. In our evolutionary model, we will assume that each individual knows only his own payoff function when choosing an action, and never randomizes. Hence $\Sigma = \prod_{i=1}^n \Sigma_i$ where Σ_i is the set of maps $\sigma^i : \Theta_i \rightarrow A_i$.

Whenever μ is fixed, we will view Σ as a topological space of random variables on the measure space (Θ, μ) with the topology τ_μ of convergence in measure.⁵ Given a profile of strategies $\sigma \in \Sigma$, the *utility* of player i is the random variable $\theta^i(\sigma(\theta))$.⁶ The fitness of player i is the random variable $\pi^i(\sigma(\theta))$. The *outcome* of play, denoted $x_\mu(\sigma)$, is the distribution of σ , which is an element of Δ .

In the model we have described, a player cares only about the distribution over opponents' actions, not the opponents' types. We will therefore simplify notation by defining best-replies as functions of outcomes, rather than strategy profiles. The pure action best-reply correspondence for a given game θ is

$$B_\theta(x) := \prod_{i=1}^n \operatorname{argmax}_{a \in A_i} \theta^i(a, x^{-i})$$

and the pure strategy best-reply correspondence in $\Gamma(\mu)$ is

$$\beta_\mu(x) := \prod_{i=1}^n \operatorname{argmax}_{\sigma_i \in \Sigma_i} \mathbf{E}_\mu(\theta^i(\sigma_i, x^{-i}))$$

In the above notation, \mathbf{E}_μ denotes expectation with respect to the measure μ , and $\theta^i(\sigma^i, x^{-i})$ is the random utility to player i when using strategy σ^i against the opponents' distribution of play x^{-i} .

We assume that the aggregate distribution of play can be described by an equilibrium in $\Gamma(\mu)$. A pure-strategy equilibrium in $\Gamma(\mu)$ is a profile σ of pure strategies with distribution x such that for each i , $\sigma^i \in \beta_\mu(x)$.

Because we have restricted attention to non-atomic distributions, the results of Milgrom and Weber (1985) and Radner and Rosenthal (1982) imply that the restriction to pure-strategies entails no loss of generality. In particular, for any $\mu \in \Theta$, there is at least one equilibrium of $\Gamma(\mu)$ in pure strategies, and for any "mixed" equilibrium there is a "purification," i.e. a pure strategy equilibrium which is equivalent in all relevant respects.

We conclude this section with the following useful lemma.

⁵Convergence in measure is the appropriate topology for our purposes because strategies are interpreted as aggregate action profiles. Thus, two strategies are "close" only if they are point-wise close for a large fraction of the population.

⁶Thus, $\Gamma(\mu)$ has private values. An interesting extension is the case in which matched individuals obtain additional information about one another's' preferences. This would capture a world in which traits such as kindness, self-interest, and ruthlessness are (perhaps imperfectly) perceptible.

Lemma 1 1. For every $\theta \in \Theta$, and $x \in \Delta$, there is a $\delta_\theta > 0$ such that $\|\hat{x} - x\| < \delta_\theta \Rightarrow B_\theta(\hat{x}) \subset B_\theta(x)$.

2. If $\mu_1 \ll \mu_2$, then $\beta_{\mu_2} \subset \beta_{\mu_1}$

3. β_μ is upper hemi-continuous.

Proof: The first point follows immediately from the continuity of the utility functions θ . The second is due to the fact that for any given opposing distribution x , two best-replies can differ on a set of measure zero, but nowhere else.

To prove the third, it suffices to show that for every $\varepsilon > 0$, there exists a $\delta > 0$ such that if $\|\hat{x} - x\| < \delta$ and $\sigma \in \beta_\mu(\hat{x})$ then $\mu(\{\theta : \sigma(\theta) \in B_\theta(x)\}) > 1 - \varepsilon$. Let X be a subset of Θ of μ -measure at least $1 - \varepsilon$ and let $\delta = \inf_X \delta_\theta$ where δ_θ is as defined in the first part of this Lemma. Then if $\|\hat{x} - x\| < \delta$, and if $\sigma \in \beta_{\mu^*}(\hat{x})$ then $\mu^*(\{\theta : \sigma(\theta) \in B_\theta(x)\}) \geq \mu^*(X) \geq (1 - \varepsilon)$. ■

3 Mutation and Selection

The central question of this paper is whether evolutionary forces, acting on the preferences in the population, will bring about distributions μ such that equilibrium play $\sigma \in \mathcal{E}(\mu)$ is in E , i.e. corresponds to a Nash equilibrium of the true game. In our analysis of this question, we do not develop an explicit model of the evolutionary process. Instead, we follow in the spirit of “static” notions of evolutionary stability such as ESS Maynard Smith (1982) (see also Weibull (1995) Chp 2 for a survey of this approach.) That is, we propose criteria which characterize “stable” sets of preference distributions, and argue that these criteria capture the important features of unmodeled evolution.

Our criterion for evolutionary stability of preferences has the usual two components: *natural selection*, the process by which unsuccessful types are replaced by successful types, and *mutation*, the process by which previously unrepresented types can enter the population. Roughly, an outcome $x \in \Delta$ is supported by stable preferences if the preferences which support the outcome are stable under natural selection, and robust to mutation. We discuss these features in the present section.

Our representation of mutations is a generalization of the representation found in traditional concepts of evolutionary stability. ESS, for example, tests the stability of an outcome by ensuring its robustness against small perturbations of the population strategy profile. Essentially this amounts to verifying that evolutionary forces will restore the original profile starting from any profile in some arbitrarily small neighborhood. The implicit idea is that players’ strategies

are subject to mutation, but we can be sure that the aggregate profile cannot move far before forces of natural selection come to dominate.

In our model, evolutionary forces operate on the distribution of preferences in the population. We therefore need to characterize the types of preference distributions which can come about as a consequence of mutation, starting from an arbitrary initial distribution μ . There are two criteria for a definition of neighborhoods of post-mutation distributions. The first is that they be sufficiently rich so as to allow all combinations of types to enter the population. The second is that they be “bounded,” capturing the idea that mutation operates slowly.

A notion of closeness which meets these criteria is the following.

Definition 1 *A neighborhood of a type distribution $\mu \in \mathcal{P}(\Theta)$ is a set U of full-support distributions such that for some $\varepsilon > 0$, $\|\mu - \nu\| < \varepsilon$ for all $\nu \in U$. If S is a set of priors, then a neighborhood of S is the union of neighborhoods of the elements of S .*

Individuals in the population choose actions which are optimal relative to their preferences. However, the objective successfulness of an action is determined by the true payoff function π . A process of natural selection operates on the types in the population, favoring those types which are most “fit” relative to the current environment. Formally, we assume that the distribution μ evolves according to the relative successfulness, evaluated according to π , of the equilibrium actions being used by individuals of various types.

Our goal is to incorporate as much of the consequences of such dynamics as possible without restricting attention to any specific process. To do this, we construct sequences of distributions, called *selection sequences* which can be considered abstract “paths” of the evolutionary process. Each step of the sequence is assumed to satisfy a standard “payoff monotonicity” property (see Weibull (1995)).

Definition 2 *A pair of distributions (μ_1, μ_2) is a **selection step** relative to σ if $\sigma \in \mathcal{E}(\mu_1)$, $\mu_2 \ll \mu_1$ and for each $i \in \{1, \dots, n\}$ and for every $X, Y \subset \Theta_i$, such that $\mu_1(X), \mu_1(Y) > 0$,*

$$\mathbf{E}_{\mu_1}(\pi^i(\sigma)|X) \begin{pmatrix} > \\ = \\ < \end{pmatrix} \mathbf{E}_{\mu_1}(\pi^i(\sigma)|Y) \Rightarrow \frac{\mu_2(X)}{\mu_1(X)} \begin{pmatrix} \geq \\ = \\ \leq \end{pmatrix} \frac{\mu_2(Y)}{\mu_1(Y)}$$

Evolution will be assumed to continue until a distribution and equilibrium are reached that are invariant, i.e. any selection step is trivial.

Definition 3 A distribution μ is **stable** with respect to outcome $x \in \Delta$ if there is some equilibrium $\sigma \in \mathcal{E}(\mu)$, whose distribution is x and such that there exist constants c_i such that $\pi^i(\sigma) = c_i$ μ -almost surely for every $i \in \{1, \dots, n\}$.

The notation $\delta(\mu)$ will represent the set of outcomes with respect to which μ is stable. If S is a set of distributions, then $\delta(S)$ is the set $\cup_{\mu \in S} \delta(\mu)$. Say that the set S is stable with respect to the set of outcomes O if $\delta(\mu) \neq \emptyset$ for each $\mu \in S$ and if $\delta(S) = O$.

In general, stability with respect to outcome x implies nothing about x . However, as long as the support of the type distribution intersects $\tilde{\pi}$, the class of preferences which are equivalent to the true preferences, stability implies that x is a Nash equilibrium.

Proposition 1 Suppose $\tilde{\pi} \cap C(\mu) \neq \emptyset$ and x is the distribution of some equilibrium $\sigma \in \mathcal{E}(\mu)$. Then μ is stable with respect to x if and only if x is a Nash equilibrium of G

To reach a stable preference distribution, selection may have to operate for more than a finite number of steps. We define a selection sequence to be any infinite sequence of distributions which satisfy the payoff-monotonicity property.

Definition 4 A sequence of pairs $\{(\mu_k, \sigma_k)\}$ is a **selection sequence** if for every k , (μ_k, μ_{k+1}) is a selection step relative to σ_k .

Say that a selection sequence $\{(\mu_k, \sigma_k)\}$ converges if there is a distribution μ^* such that $\mu_k \rightarrow \mu^*$ in norm. Because we have assumed little more than payoff monotonicity in selection sequences, there is no guarantee that a limit point of a selection sequence may itself be stable. In other words, it is possible that the selection will continue to occur once the limit is reached. We could impose assumptions which rule out such ill-behaved selection sequences (for example, focusing on a particular class of dynamic processes which satisfy a continuity assumption). Instead, in the interest of maintaining as much generality as possible, we will simply focus on the stable limit points, implicitly assuming that selection continues to occur following convergence to an unstable limit point.⁷ Let $\mathcal{L}(\mu)$ be the set of stable limit points of selection sequences starting with μ . We will later show (Proposition 4) that $\mathcal{L}(\mu)$ is non-empty under quite general circumstances.

⁷One may still be concerned that the limit point μ^* might be stable with respect to some x , yet play along the sequence never approaches x . Proposition 4 shows that mild assumptions guarantee that μ^* is stable with respect to at least one accumulation point of the sequence of play.

4 Stability

In this section we introduce the first of two stability definitions. We say that a set of outcomes O is supported by stable preferences if there is a set S of preference distributions which are stable with respect to that set of outcomes, and a neighborhood U of S from which selection must return to S . Without imposing any further restrictions on selection sequences we show that every game has a set of outcomes that is supported by stable preferences and every such set is contained in the set of Nash equilibria.

Definition 5 *A set $O \subset \Delta$ of outcomes is supported by stable preferences if it is a minimal non-empty closed set with the following property. There exists a subset $S \subset \mathcal{P}(\Theta)$ which is stable with respect to O and a neighborhood U of S such that $\emptyset \neq \mathcal{L}(\nu) \subset S$ for all $\nu \in U$.*

We begin with an existence result.

Theorem 1 *Every game has at least one set of outcomes that is supported by stable preferences.*

The proof of this theorem relies on some properties of selection sequences which we now establish.

Proposition 2 *Suppose $\mu_0(\tilde{\pi}) > 0$. Then every selection sequence beginning with μ_0 converges to a limit distribution μ^* satisfying $\mu^*(\tilde{\pi}) > 0$.*

Proof: Let $\{(\mu_k, \sigma_k)\}$ be a selection sequence beginning with μ_0 . For every k and every equilibrium of $\Gamma(\mu_k)$, the types in $\tilde{\pi}$ play actions that maximize π . Therefore, by the definition of a selection step, the sequence $\mu_k(\tilde{\pi})$ is weakly increasing. Since $\mu_k(\tilde{\pi}) \in (0, 1]$ the sequence must converge, implying that

$$z_k := \frac{\mu_k(\tilde{\pi})}{\mu_{k-1}(\tilde{\pi})} \rightarrow 1$$

To show that μ_k converges, we will show that it converges in norm, i.e. that $\|\mu_k - \mu_{k-1}\| \rightarrow 0$. For this it is sufficient to show $\|(\mu_k - \mu_{k-1})^+\| \rightarrow 0$. Let X be the support of $(\mu_k - \mu_{k-1})^+$. By the definition of a selection step

$$\frac{\mu_k(X)}{\mu_{k-1}(X)} \leq z_k$$

implying

$$\begin{aligned} \|(\mu_k - \mu_{k-1})\| &\equiv \mu_k(X) - \mu_{k-1}(X) \\ &\leq (z_k - 1)\mu_{k-1}(X) \\ &\leq z_k - 1 \end{aligned}$$

and we have shown that the right-hand side converges to zero. ■

Proposition 3 *Suppose $\mu_0(\tilde{\pi}) > 0$. Then there exists a selection sequence beginning with μ_0 whose limit is stable.*

Proof: This is a consequence of Proposition 4 which will be proven in Section 5. ■

Proof of Theorem 1 Let G be a game with payoff function π . Every Nash equilibrium of G can be supported as an equilibrium of $\Gamma(\delta_{\tilde{\pi}})$ where $\delta_{\tilde{\pi}}$ is any prior concentrated on $\tilde{\pi}$. By Proposition 1, these distributions are stable. Therefore the set S of priors μ that are stable with respect to the set E and for which $\mu(\tilde{\pi}) > 0$ is non-empty. Furthermore, for every neighboring prior ν $\nu(\tilde{\pi}) > 0$ and Propositions 3 and 1 imply $\mathcal{L}(\nu) \neq \emptyset$ and $\mathcal{L} \subset S$.

Thus the closed set E satisfies the criteria in the definition, and by the usual Zorn's lemma argument (for example, see Kohlberg and Mertens (1986, Proposition 1) and Kalai and Samet (1984, Theorem 1)) there is a minimal closed subset of E which does as well. ■

Only Nash equilibria can be stable.

Theorem 2 *A set O is supported by stable preferences only if $O \subset E$.*

Proof: Let O be a set of outcomes that is supported by stable preferences. Then there is a set S of priors which are stable with respect to O and a neighborhood U_μ of each $\mu \in S$ such that for every $\nu \in U_\mu$,

$$\mathcal{L}(\nu) \subset S \tag{1}$$

Every such neighborhood contains a $\nu \in F \equiv \{\hat{\nu} : \hat{\nu}(\tilde{\pi}) > 0\}$. By Propositions 2 and 3 we have $\emptyset \neq \mathcal{L}(\nu) \subset F$. Therefore $F \cap S \neq \emptyset$.

Proposition 1 implies $\delta(F \cap S) \subset E$, and by definition $\delta(F \cap S) \subset O$. Let $Q \equiv E \cap O$. Q is closed because both E and O are, and because $F \cap S \neq \emptyset$, Proposition 1 implies $Q \neq \emptyset$.

Suppose $\delta(F \cap S) = Q$. Then Q satisfies the criteria for supported by stable preferences using the neighborhood $\cup_{\mu \in F \cap S} U_\mu$ of $F \cap S$.

On the contrary, if $x \in Q \setminus \delta(F \cap S)$, then there is a $\mu \in S \setminus F$ which is stable with respect to x . Let ι be any distribution concentrated on $\tilde{\pi}$. Since x is a Nash equilibrium, all actions with positive probability are best-replies under $\tilde{\pi}$ so that any strategy which has distribution x in $\Gamma(\iota)$ is an equilibrium of $\Gamma(\iota)$ (such strategies exist because distributions are assumed atomless).

Lemma 2 implies that for any $s \in (0, 1)$, the probability

$$\nu := (1 - s)\mu + s\iota$$

has an equilibrium γ whose distribution is x . For s small enough ν is inside U_μ . Let U_ν be a neighborhood of ν contained in U_μ with radius no greater than s . Thus $\emptyset \neq \mathcal{L}(U_\nu) \subset \mathcal{L}(U_\mu) \subset S \cap F$ by Lemma 3 and by (1).

Thus, $S \cap F$ together with all distributions constructed in this manner constitute a stable set of priors that support Q . Finally, since O is supported by stable preferences it is minimal, hence $O = Q \subset E$. ■

5 Perfect Equilibrium

In the previous section we established that under very general conditions, outcomes which are supported by stable preferences must be Nash equilibria. In this section we show that mild additional assumptions ensure that only trembling-hand perfect equilibria can be supported by stable preferences.

We first impose some weak regularity conditions on selection sequences. The first condition, which we call *bounded death rates* ensures that selection does not stop “too early.” In particular, the rate at which successful types grow is asymptotically bounded below. The second condition, *finite death rates* prevents types from becoming completely extinct in one step. A feature of this assumption which plays an important role in this section is that types for whom a given strategy is dominant can never be completely eliminated.

As a final modification of this section, we pay more attention to the sequence of play along an evolutionary path. Suppose μ^* is the limit of a selection sequence $\{(\mu_k, \sigma_k)\}$ within some stable set of distributions. We restrict the equilibria of μ^* to those which are approximated asymptotically by the sequence of equilibria σ_k . In doing so, we are implicitly assuming a sort of continuity in the evolution of equilibrium play. We do not use μ^* as support for some outcome x for which x is not an accumulation point of the distribution of play along the sequence.

Given a selection sequence $\{(\mu_k, \sigma_k)\}$, let x_k be the distribution of play in period k , and let $W_k = \{\theta : \sigma_k(\theta) \notin B_\pi(x_k)\}$. This is the set of types which are not maximizing fitness in period k .

Definition 6 A selection sequence $\{(\mu_k, \sigma_k)\}$ has **bounded death rates** if

$$\limsup_{k \rightarrow \infty} \frac{\mu_{k+1}(W_k)}{\mu_k(W_k)} < 1$$

Definition 7 A selection sequence $\{(\mu_k, \sigma_k)\}$ has **finite death rates** if $\mu_k \ll \mu_{k+1}$ for every k .

These assumptions are weak and are satisfied for example by the replicator dynamics. The following characterizes the behavior of such selection sequences.

Proposition 4 Assume $\mu_0(\tilde{\pi}) > 0$. If $\{(\mu_k, \sigma_k)\}$ is a selection sequence beginning with μ_0 with finite and bounded death rates, then there exists a stable μ^* such that $\mu_k \rightarrow \mu^*$. Moreover, every accumulation point of x_k is a perfect equilibrium of G which is the distribution of some equilibrium of $\Gamma(\mu^*)$.

Proof: For every k we have

$$\frac{\mu_k(\tilde{\pi})}{\mu_{k-1}(\tilde{\pi})} = \frac{\mu_k(W_{k-1}^c)}{\mu_{k-1}(W_{k-1}^c)}$$

Subtracting 1 from both sides

$$\frac{\mu_k(\tilde{\pi}) - \mu_{k-1}(\tilde{\pi})}{\mu_{k-1}(\tilde{\pi})} = \frac{\mu_k(W_{k-1}^c) - \mu_{k-1}(W_{k-1}^c)}{\mu_{k-1}(W_{k-1}^c)}$$

Now

$$\begin{aligned} \mu_k(W_{k-1}^c) - \mu_{k-1}(W_{k-1}^c) &= \mu_{k-1}(W_{k-1}) - \mu_k(W_{k-1}) \\ &= \left(1 - \frac{\mu_k(W_{k-1})}{\mu_{k-1}(W_{k-1})}\right) \mu_{k-1}(W_{k-1}) \end{aligned}$$

Under the assumption of bounded death rates, there exists $\varepsilon > 0$ and \bar{k} such that for all $k > \bar{k}$, the right-hand side is greater than $(1 - \varepsilon)\mu_{k-1}(W_{k-1})$. Therefore, for all such k ,

$$\frac{\mu_k(\tilde{\pi}) - \mu_{k-1}(\tilde{\pi})}{\mu_{k-1}(\tilde{\pi})} > \frac{(1 - \varepsilon)\mu_{k-1}(W_{k-1})}{1 - \mu_{k-1}(W_{k-1})} \geq 0$$

Taking limits as k goes to infinity, Proposition 2 implies that $\mu_k \rightarrow \mu^*$, hence the left-hand side converges to 0. We therefore conclude

$$\lim_{k \rightarrow \infty} \mu_k(W_k) = 0$$

Now let x be an accumulation point of the sequence x_k of play, and consider a subsequence $\{(\mu_l, \sigma_l)\}$ of $\{(\mu_k, \sigma_k)\}$ whose play is converging to x .

We have shown that $\mu_l(W_l) \rightarrow 0$. Thus, for every $\varepsilon > 0$ there exists a $\bar{l}(\varepsilon)$ such that $l > \bar{l}(\varepsilon)$ implies $\mu_l(W_l) < \varepsilon$. Furthermore, the assumption of finite death rates implies x_l is completely mixed for every l .⁸ In other words, x_l is a ε -perfect equilibrium distribution for every $l > \bar{l}(\varepsilon)$.

We can thus choose any sequence $\varepsilon_z \rightarrow 0$, and the corresponding sequence $x_{\bar{l}(\varepsilon_n)}$ is a sequence of ε_n -perfect equilibria. Therefore x is a perfect equilibrium outcome of G . We wish to show that there is an equilibrium of $\Gamma(\mu^*)$ whose distribution is x .

The upper hemi-continuity of β_{μ^*} implies that for every neighborhood U of $\beta_{\mu^*}(x)$ there is a \bar{l} such that $l > \bar{l}$ implies $\beta_{\mu^*}(x_l) \subset U$. Now $\mu^* \ll \mu_l$ follows from norm convergence, hence by Lemma 1 we have $\beta_{\mu_l} \subset \beta_{\mu^*}$, and therefore $\sigma_k \in \beta_{\mu_k}(x_k) \subset U$. Thus all τ_{μ^*} -accumulation points of σ_k are in $\beta_{\mu^*}(x)$. Let σ^* be any one of them. Convergence in mean implies convergence in distribution, hence x must be the distribution of σ^* . We conclude that σ^* is an equilibrium of $\Gamma(\mu^*)$ with distribution x . ■

This result motivates an alternative definition of support by stable preferences. Under the previous definition, the stable set of outcomes must include all stable equilibria of elements of the stable set of preferences. Under this alternative definition, we include only those equilibria that are accumulation points of play along the selection sequence. Proposition 4 makes it clear that such a requirement will imply a refinement of Definition 5.

For the remainder, unless otherwise noted, we will restrict attention to selection sequences with bounded, finite death rates. We now define a limit point of a selection sequence $\{(\mu_k, \sigma_k)\}$ to be a pair $(\mu^*, \sigma^*) \in \mathcal{P}(\Theta) \times \Sigma$ where μ^* is the norm limit of μ_k and σ^* is a τ_{μ^*} accumulation point of σ_k . In general, while there can be only one limit of μ_k , there may be multiple accumulation points of σ_k , hence a selection sequence can have more than one limit. Let $\mathcal{L}^*(\mu_0)$ denote the set of limits of selection sequences starting with μ_0 .

Some additional notation will come in handy. An *refinement* is a correspondence $\rho : S \rightarrow \Sigma$ whose domain is a subset S of $\mathcal{P}(\Theta)$, satisfying $\rho(\mu) \subset \mathcal{E}(\mu)$.

We will say that a pair (μ, σ) is stable with respect to outcome x if σ is an equilibrium of $\Gamma(\mu)$ with distribution x and there is a fitness profile $c = (c_i)_{i=1}^n$

⁸The easiest way to see this is to note that any full-support distribution must put positive mass on the set of types for whom action a is strongly dominant, for every a .

such that $\pi(\sigma) = c$ μ -almost surely. Represent by $\delta(\mu, \sigma)$ the set of outcomes with respect to which (μ, σ) is stable. Say that a refinement ρ is stable with respect to the set of outcomes $O \subset \Delta$ if $\delta(\varphi) \neq \emptyset$ for each $\varphi \in \text{graph}\rho$ and $\cup_{\varphi \in \text{graph}\rho} \delta(\varphi) = O$.

Definition 8 *A set $O \subset \Delta$ of outcomes is supported by stable preferences if it is a minimal non-empty closed set with the following property. There exists a $S \subset \mathcal{P}(\Theta)$ and a refinement on S which is stable with respect to O and if for each $\mu \in S$ there is a neighborhood U of μ such that $\emptyset \neq \mathcal{L}^*(\nu) \subset \text{graph}\rho$ for each $\nu \in U$, $\nu \neq \mu$.*

With Proposition 4 in hand, we can mimic the proofs of Theorems 1 and 2 to establish the following.

Theorem 3 *Every game G has a set of outcomes that is supported by stable preferences. A set of outcomes is supported by stable preferences only if it consists of perfect equilibrium distributions of G .*

6 Conclusion

We conclude with a summary of the advantages of our approach to evolutionary equilibrium foundations.

One of the main goals of evolutionary game theory has been to provide a foundation for Nash equilibrium and perhaps argue in favor of some of its refinements. Models based on evolution of *strategies* such as ESS and related dynamic models have been partially successful. Generally speaking, outcomes that satisfy these types of evolutionary stability criteria must be Nash equilibria, and in many cases must satisfy refinements incorporating backward and forward induction ideas. A major drawback of this approach, however, has been that evolutionarily stable strategies often fail to exist.

We have based our solution concept on a model of strategic interaction among rational agents whose *preferences* are subject to evolutionary forces. We find this an appealing alternative *prima facie* as a model of economic behavior. The results of this paper show that the model has further advantages. First, every game has at least one outcome that is supported by stable preferences (Theorem 1.) Additionally, we preserve the standard “only if Nash” result (Theorem 2) and show support for a traditional equilibrium refinement (Theorem 5). These results are obtained by imposing a minimum of structure on the model of natural selection.

Finally, our model suggests a natural solution to one of the fundamental conceptual difficulties of equilibrium theory: the interpretation of mixed strategies. In our model, individuals never randomize. Mixed outcomes only appear random to an observer *outside* the model who has no information about the exact preferences of the individuals playing G . Our model thus demonstrates how evolution of preferences leads to purification of mixed equilibria.

References

- BINMORE, K., L. SAMUELSON, AND R. VAUGHN (1992): "Musical Chairs: Modeling Noisy Evolution," *Games and Economic Behavior*, 6, 100–200.
- BINMORE, K. G., AND L. SAMUELSON (1997): "Muddling Through: Noisy Equilibrium Selection," *Journal of Economic Theory*, 71, 235–265.
- ELY, J. C. (1996): "Local Conventions," mimeo, Northwestern University.
- KALAI, E., AND D. SAMET (1984): "Persistent Equilibria in Strategic Games," *International Journal of Game Theory*, 13, 129–144.
- KANDORI, M., G. MAILATH, AND R. ROB (1993): "Learning, Mutation, and Long Run Equilibria in Games," *Econometrica*, 65, 383–414.
- KOHLBERG, E., AND J.-F. MERTENS (1986): "On the Strategic Stability of Equilibria," *Econometrica*, 54, 1003–1037.
- MAYNARD SMITH, J. (1982): *Evolution and the Theory of Games*. Cambridge University Press, Cambridge.
- MILGROM, P. R., AND R. J. WEBER (1985): "Distributional Strategies for Games with Incomplete Information," *Mathematics of Operations Research*, 10, 619–632.
- RADNER, R., AND R. W. ROSENTHAL (1982): "Private Information and Pure Strategy Equilibria," *Mathematics of Operations Research*, 7, 401–409.
- ROBSON, A. J. (1996a): "A Biological Basis for Expected and Non-Expected Utility," *Journal of Economic Theory*, 68, 397–424.
- (1996b): "The Evolution of Attitudes to Risk," *Games and Economic Behavior*, 14, 190–207.
- ROBSON, A. J., AND F. VEGA-REDONDO (1996): "Efficient Equilibrium Selection in Evolutionary Games with Random Matching," *Journal of Economic Theory*, 70, 43–64.
- SANDHOLM, W. H., AND A. PAUZNER (1997): "Evolution, Population Growth, and History Dependence," mimeo, Northwestern University.
- TAYLOR, P. D., AND L. B. JONKER (1978): "Evolutionary Stable Strategies and Game Dynamics," *Mathematical Biosciences*, 40, 145–56.

TO, T. (1995): "Risk and Evolution," mimeo, St. Andrews University.

WEIBULL, J. W. (1995): *Evolutionary Game Theory*. MIT Press, Cambridge.

YOUNG, H. P. (1993): "The Evolution of Conventions," *Econometrica*, 61, 57–84.

A Omitted Proofs

Lemma 2 *Suppose μ, ν are distributions each of which has an equilibrium whose distribution is x . Then for any $s \in (0, 1)$, the distribution $s\mu + (1 - s)\nu$ has an equilibrium whose distribution is x .*

Proof: Let σ and γ be the equilibria of μ and ν respectively. Consider the behavior strategy ⁹ $\tilde{\sigma}(\theta) = s\sigma(\theta) + (1 - s)\gamma(\theta)$. This is an equilibrium of $\Gamma(s\mu + (1 - s)\nu)$ in behavior strategies because each type is randomizing over best-replies. Clearly its distribution is x . Milgrom and Weber (1985) prove that any such equilibrium has a *purification*, i.e. a pure strategy equilibrium with the same distribution. ■

⁹A behavior strategy is a map $\tilde{\sigma} : \Theta \rightarrow \Delta$ specifying a mixed-strategy distribution for each type.