



Munich Personal RePEc Archive

# **Stability and trembles in extensive-form games**

Heller, Yuval

University of Oxford

9 July 2013

Online at <https://mpra.ub.uni-muenchen.de/50350/>

MPRA Paper No. 50350, posted 02 Oct 2013 18:38 UTC

# Stability and Trembles in Extensive-form Games

Yuval Heller\* (October 2, 2013)

Address: Nuffield College and Department of Economics, New Road, Oxford, OX1 1NF, United Kingdom. Email: yuval26@gmail.com or yuval.heller@economics.ox.ac.uk.

## Abstract

A leading solution concept in the evolutionary study of extensive-form games is Selten's (1983) [16] notion of limit ESS. This note demonstrates that a limit ESS does not imply neutral stability, and that it may be dynamically unstable (almost any small perturbation takes the population away). These problems arise due to an implicit assumption that “mutants” are arbitrarily rare relative to “trembling” incumbents. Finally, I present a novel definition that solves this issue and has appealing properties.

KEYWORDS: Limit ESS, evolutionary stability, extensive-form games. JEL Classification: C73.

## 1 Introduction

In a seminal paper, Maynard-Smith & Price [12] defined an evolutionarily stable strategy (ESS) as a Nash equilibrium that is a *strictly* better reply against other best-reply strategies. It was extended in [11] to the weaker notion of a neutrally stable strategy (NSS) that is a *weakly* better reply against other best-reply strategies. The motivation for these notions is that a stable strategy, if adopted by a population of players, cannot be invaded by any alternative strategy that is initially rare. This is formalized in [4, 17], in which it is shown that any NSS is *Lyapunov stable* in the replicator dynamics: no small change in the population composition can take it away from the state in which everyone follows the NSS, and any ESS is *asymptotically stable*: any sufficiently small change results in a movement back toward the ESS.

Extensive-form games rarely admit an ESS due to the existence of “equivalent” strategies that differ only off the equilibrium path. Selten [16] relaxes this notion by requiring evolutionary stability in a converging sequence of perturbed games in which the players may infrequently

---

\*I would like to express my deep gratitude to Erik Mohlin and Thomas Norman for many useful discussions.

“tremble” and play different actions by mistake (see Section 2 for the formal definitions). Selten’s solution concept, called *limit ESS*, is a central notion in the evolutionary study of games with more than one stage, and it has been applied to various interactions in the economics and biology literature (see, e.g., [3, 6, 8, 10]).<sup>1</sup>

van Damme [18, Cor. 9.8.6.iii] has proved that the notion of limit ESS refines neutral stability.<sup>2</sup> However, Van Damme’s proof contains a mistake (as explained in footnote 4). Section 3 presents a simple game that admits a limit ESS that: (1) is not neutrally stable, and (2) is dynamically unstable in a strong sense: almost any nearby initial state takes the population away from the limit ESS. Section 4 shows that the reason for this is the implicit assumption in the notion of limit ESS that mutants are arbitrarily rare relative to the “trembling” incumbents. Finally, I present a novel definition, which I call a *uniform limit ESS*, that only assumes that the mutants are sufficiently rare relative to the non-trembling incumbents. I show that this new notion refines limit ESS, implies neutral stability, and has appealing properties.

## 2 Definitions

Let  $\Gamma$  be a symmetric two-player extensive-form game (a formal detailed definition is given in the appendix). Let index  $i \in \{1, 2\}$  denote one of the players, and let  $-i$  denote the other player. Let  $U_i$  be the set of information sets of player  $i$ . For each such information set  $u \in U_i$ , let  $C_u$  be the set of *choices* (or *actions*) in information set  $u$ . The game is endowed with a symmetry function  $T$  that maps each choice  $c$  of player  $i$  to the symmetric choice  $c^T$  of player  $-i$ . Let  $B_i$  denote the set of *behavior strategies* of player  $i$  (a mapping that assigns a probability distribution over the set of choices at each information set of player  $i$ ). Let  $R_i(b_1, b_2)$  be the expected payoff to player  $i$  when each player  $i$  plays strategy  $b_i \in B_i$ . Given strategy  $b \in B_1$ , let  $b^T$  denote the symmetric strategy of player 2. The symmetry between the strategies implies that  $R_1(b, b^T) = R_2(b, b^T)$ .

An evolutionarily (neutrally) stable strategy (abbreviated ESS, NSS) is a strategy that satisfies two conditions: (1) it is a best reply to itself (i.e., a symmetric Nash equilibrium), and (2) it achieves a strictly (weakly) better payoff against any other best-reply strategy. Formally:<sup>3</sup>

**Definition 1.** [11, 12] Strategy  $b \in B_1$  is an ESS (NSS) if for every  $\tilde{b} \in B_1$  ( $\tilde{b} \neq b$ ):

1.  $R_1(b, b^T) \geq R_1(\tilde{b}, b^T)$ ; and

<sup>1</sup> The notion is also central in the study of asymmetric one-shot games that are played by a population in which each agent is randomly assigned one of the roles in the game (see, e.g., [7, 14]).

<sup>2</sup> Accordingly, several papers state that “it is well known that a limit ESS is an NSS (e.g., Bhaskar [1, page 274] and Bhaskar [2, page 115]).

<sup>3</sup> Part of the literature calls it “direct-ESS,” and uses the name “ESS” only for mixed strategies.

2. if  $R_1(b, b^T) = R_1(\tilde{b}, b^T)$ , then  $R_1(b, \tilde{b}^T) > R_1(\tilde{b}, \tilde{b}^T)$  ( $R_1(b, \tilde{b}^T) \geq R_1(\tilde{b}, \tilde{b}^T)$ ).

Bomze & Weibull [4] showed that any NSS is *Lyapunov stable* in the replicator dynamics: populations starting close enough to the NSS remain close forever (though a sequence of small perturbations may take the population away). Taylor & Jonker [17] showed that ESS satisfies the stronger notion of *asymptotic stability*: populations starting close enough to the ESS eventually converge to it (see extensions to other payoff-monotonic dynamics in [5, 15]).

Extensive-form games rarely admit an ESS due to the existence of “equivalent” strategies that differ only off the equilibrium path. Selten [16] relaxes this notion by requiring evolutionary stability only in a converging sequence of perturbed games (but not necessarily in the unperturbed game). Formally:

**Definition 2.** [16] A perturbation of a symmetric two-player extensive-form game  $\Gamma$  is a mapping  $\eta$  from the set of choices into the reals such that: (1) for each choice  $c$  the following hold:  $\eta(c) \geq 0$  and  $\eta_c = \eta_{c^T}$ ; and (2) for each information set  $u$ :  $\sum_{c \in C_u} \eta(c) < 1$ .

The perturbed game  $(\Gamma, \eta)$  has the same structure as  $\Gamma$  except that strategy  $b$  is admissible only if  $b_u(c) \geq \eta_c$  for all  $u$  and  $c$ . Let  $B_i(\eta)$  denote the set of all such admissible strategies of player  $i$ . A limit ESS is the limit point of the ESS of a converging sequence of perturbed games.

**Definition 3.** [16] Strategy  $b \in B_1$  is a *limit ESS* if there exists a sequence  $(\eta^k, b^k)_{k \in \mathbb{N}}$  with  $\eta^k \rightarrow 0$  and  $b^k \rightarrow b$  when  $k \rightarrow \infty$  such that  $b^k \in B_1(\eta^k)$  is an ESS of the game  $(\Gamma, \eta^k)$ .

Note that the special case of  $\eta \equiv 0$  is not excluded; hence, every ESS is a limit ESS.

### 3 Example: Limit ESS That is not Neutrally Stable

This section presents a limit ESS that (1) is not neutrally stable, and (2) is dynamically unstable in a strong sense.

Consider the following one-shot symmetric two-player game in which each player has to simultaneously choose either  $c_1, c_2, c_3$ , and the payoff matrix is given by Table 1. With a slight abuse of notation, let  $c_i$  denote the strategy that assigns probability one to choice  $c_i$ . Observe, first, that strategy  $c_1$  is a limit ESS. Let the sequence  $(\eta^k, b^k)$  be defined as follows for each  $k \geq 4$ :  $\eta^k(c_1) = \eta^k(c_2) = \eta^k(c_3) = \frac{1}{k}$ , and  $b^k(c_1) = 1 - \frac{2}{k}$ ,  $b^k(c_2) = b^k(c_3) = \frac{1}{k}$ . Observe that each  $b^k$  is an ESS of the perturbed game  $(\Gamma, \eta^k)$  with  $\eta^k \rightarrow 0$  and  $b^k \rightarrow c_1$  when  $k \rightarrow \infty$ .

Next, observe that strategy  $c_1$  is not an NSS, because  $R_1(c_2, c_1) = R_1(c_1, c_1)$  and  $R_1(c_2, c_2) > R_1(c_1, c_2)$ . Moreover, strategy  $c_2$  is dynamically unstable in the replicator dynamics: any initial state that assigns a positive mass to  $c_2$  takes the population in the long run to assign mass one

Tab. 1: Payoff Matrix of a Symmetric Two-Player Game

|       | $c_1$ | $c_2$ | $c_3$ |
|-------|-------|-------|-------|
| $c_1$ | 2 2   | 1 2   | 4 0   |
| $c_2$ | 2 1   | 3 3   | 1 0   |
| $c_3$ | 0 4   | 0 1   | 0 0   |

to  $c_2$ . The reason for this is that strategy  $c_3$  is strictly dominated and its frequency converges to zero; as soon as the frequency of  $c_3$  is sufficiently small, strategy  $c_2$  achieves a strictly higher payoff than all other strategies in all the remaining stages.<sup>4</sup>

## 4 Uniform Limit ESS

In this section I reformulate the definition of limit ESS to highlight what leads to the counter-intuitive implication demonstrated in the example, and I present a refinement that deals with this issue.

It is well known (see, e.g., [19, Prop. 2.1 and 2.5]) that a strategy is an ESS iff it outperforms any other strategy in a mixed population, provided that the share of the other strategy is sufficiently small. Formally:<sup>5</sup>

**Fact 1.** *Strategy  $b \in B_1$  is an ESS in a symmetric two-player game  $\Gamma$  iff there exists some  $\bar{\epsilon} \in (0, 1)$  such that for every strategy  $\tilde{b} \in B_1$  ( $\tilde{b} \neq b$ ) and every  $\epsilon \in (0, \bar{\epsilon})$ :*

$$r_1(b, \epsilon \cdot \tilde{b}^T + (1 - \epsilon) \cdot b^T) > r_1(\tilde{b}, \epsilon \cdot \tilde{b}^T + (1 - \epsilon) \cdot b^T),$$

<sup>4</sup> The example also shows the mistake in [18, Cor. 9.8.6.iii]’s proof that every limit ESS is neutrally stable. The penultimate sentence before Cor. 9.8.6 in page 253 of [18] states that if a strategy is a limit of ESSs in a converging sequence of perturbed normal-form games, then by taking the limit, it follows that the strategy is neutrally stable in the unperturbed game. That is, it says that if for each  $k$  (i)  $R_1(b^k, b^k) \geq R_1(\tilde{b}^k, b^k)$  and (ii)  $R_1(b^k, b^k) = R_1(\tilde{b}^k, b^k) \Rightarrow R_1(b^k, \tilde{b}^k) > R_1(\tilde{b}^k, \tilde{b}^k)$ , then  $c_1 = \lim_{k \rightarrow \infty} b^k$  is neutrally stable. However, this claim is false. For each  $k$ , let  $\tilde{b}^k(c_1) = \tilde{b}^k(c_3) = \frac{1}{k}$ , and  $\tilde{b}^k(c_2) = 1 - \frac{2}{k}$ , with the limit  $c_2 = \lim_{k \rightarrow \infty} \tilde{b}^k$ . Observe that for each  $k$ ,  $R_1(b^k, b^k) > R_1(\tilde{b}^k, b^k)$ , while in the limit  $R_1(c_1, c_1) = R_1(c_2, c_1)$ . Together with the observation that  $R_1(b^k, \tilde{b}^k) < R_1(\tilde{b}^k, \tilde{b}^k)$  for each  $k$  (as well as in the limit), it implies that  $c_1$  is not neutrally stable.

<sup>5</sup> The equivalence of the definitions hold for finite games. In games with an infinite number of actions (or non-linear payoffs), the stability property described in Fact 1, which assumes the existence of a uniform barrier ( $\bar{\epsilon}$ ), may be strictly stronger than Definition 1.

where  $\epsilon \cdot \tilde{b}^T + (1 - \epsilon) \cdot b^T \in B_2$  is the strategy that follows  $b^T$  with probability  $1 - \epsilon$  and follows  $\tilde{b}^T$  with the remaining probability  $\epsilon$ . The strategy is an NSS if the strict inequality is replaced by a weak inequality.

This allows us to reformulate the definition of a limit ESS as follows:

**Fact 2.** Strategy  $b \in B_1$  is a limit ESS if there exists a sequence  $(\eta^k, b^k)_{k \in \mathbb{N}}$  with  $\eta^k \rightarrow 0$  and  $b^k \rightarrow b$  when  $k \rightarrow \infty$  such that for each  $k$ : (1)  $b_k \in B_1(\eta^k)$ , and (2) there exists  $\bar{\epsilon}_k \in (0, 1)$  such that for every  $\tilde{b}_k \in B_1(\eta^k)$  ( $\tilde{b}_k \neq b_k$ ) and every  $\epsilon \in (0, \bar{\epsilon}_k)$ :

$$r_1(b_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T) > r_1(\tilde{b}_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T).$$

The order of quantifiers in the above definition implies that the share of the “mutants” who follow the different strategy can be arbitrarily low relative to the frequency of “trembling” incumbents. This is what allows strategy  $c_1$  to be a limit ESS in the above example (as the mutants’ loss against trembling incumbents outweighs their gain against other mutants).

I now present an alternative notion (which I call uniform limit ESS) that uses a uniform bound to the frequency of the mutants to capture the idea that mutants are rare relative to the incumbents, but not relative to the “trembling” incumbents. Formally:

**Definition 4.** Strategy  $b \in B_1$  is a *uniform limit ESS* if there exist  $\bar{\epsilon} \in (0, 1)$  and a sequence  $(\eta^k, b^k)_{k \in \mathbb{N}}$  with  $\eta^k \rightarrow 0$  and  $b^k \rightarrow b$  when  $k \rightarrow \infty$  such that for each  $k$ : (1)  $b_k \in B_1(\eta^k)$ , and (2) for every  $\tilde{b}_k \in B_1(\eta^k)$  ( $\tilde{b}_k \neq b_k$ ) and every  $\epsilon \in (0, \bar{\epsilon})$ :

$$r_1(b_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T) > r_1(\tilde{b}_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T).$$

It is immediate that any uniform limit ESS is a limit ESS. The following proposition shows that any uniform limit ESS is an NSS.

**Proposition 1.** *Any uniform limit ESS is an NSS.*

*Proof.* Let  $b \in B_1$  be a uniform limit ESS. Let  $\tilde{b} \in B_1$  be any other strategy  $\tilde{b} \neq b$ . Definition 4 implies that there exist  $\bar{\epsilon} \in (0, 1)$  and a sequence  $(\eta^k, b^k, \tilde{b}^k)_{k \in \mathbb{N}}$  with  $\eta^k \rightarrow 0$ ,  $b^k \rightarrow b$  and  $\tilde{b}^k \rightarrow \tilde{b}$  when  $k \rightarrow \infty$  such that for each  $k$  and every  $\epsilon \in (0, \bar{\epsilon})$ :

$$r_1(b_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T) > r_1(\tilde{b}_k, \epsilon \cdot \tilde{b}_k^T + (1 - \epsilon) \cdot b_k^T).$$

By continuity, the analogous weak inequality holds when  $b$  replaces  $b_k$  and  $\tilde{b}$  replaces  $\tilde{b}_k$ .  $\square$

## 5 Concluding Remarks

1. The applications of the notion of limit ESS (e.g., [3, 6, 7, 8, 10]) also satisfy the refinement of a uniform limit ESS. In this sense, the refinement is not “too strong”: it omits implausible limit ESSs like the one presented in Section 3, but it includes interesting and plausible limit ESSs in applications.
2. In [6] I presented another refinement of limit ESS, which I called *strict limit ESS*, that requires the strategy to be the limit of ESS for every converging sequence of *ubiquitous* perturbed games (which assign minimal *positive* probabilities for all choices). The motivation (which is similar to Okada [13]’s arguments for the notion of strict perfection) is that a notion stability should be robust to the specific structure of the perturbations. One can show that these two refinements are independent.
3. Prop. 1 implies that any uniform limit ESS is Lyapunov stable. I conjecture that a uniform limit ESS, which is also a strict limit ESS, satisfies a stronger notion of dynamic stability (but weaker than asymptotic stability): the share of the population who follows the uniform limit ESS strictly increases from almost any close enough initial state.

## A A Formal Detailed Definition of a Two-Player Symmetric Game

The definition is based on [16] and [18, Chapters 6 and 9], and the reader is referred to these references for interpretation and further details.

A *symmetric two-player extensive-form game* is a tuple  $\Gamma = (K, P, U, C, p, r, T)$  where:

- The *game tree*  $K$  is a finite tree with a distinguished node  $\phi$  - the root of  $K$ . Given a node in the tree  $x$ , let  $S(x)$  denote its (*immediate*) *successor*. Let  $Z$  be the endpoints of the tree (nodes with no successors), and let  $X$  be the set of nodes with successors (*decision points*). The unique sequence of nodes and branches connecting the root  $\phi$  with a node  $x$  is called the path to  $x$ . We say that  $x$  comes before  $y$  if  $x$  is on the path to  $y$  and  $x \neq y$ .
- The *player partition*  $P$  is a partition of  $X$  into 3 sets:  $P_0, P_1, P_2$ . The set  $P_i$  is the set of decision points of player  $i$ . Player 0 is the *chance* player responsible for the random moves occurring in the game.
- The *information partition*  $U$  is a pair  $(U_1, U_2)$ , where  $U_i$  is a partition of  $P_i$  (the so-called *information sets* of player  $i$ ) such that: (1) every path intersects each information set at most once, and (2) all nodes in each information set have the same number of successors.

- The *choice partition*  $C$  is a collection  $C = \{C_u | u \in U_1 \cup U_2\}$ , where  $C_u$  is a partition of  $\cup_{x \in U} S(x)$  into so-called *choices* (or *actions*) at  $u$ , such that every choice contains exactly one element of  $S(x)$  for every  $x \in U$ .
- The *probability assignment*  $p$  specifies for every  $x \in P_0$  a completely mixed probability distribution  $p_x$  on  $S(x)$ .
- The *payoff function*  $r$  is a pair  $(r_1, r_2)$ , where  $r_i : Z \rightarrow \mathbb{R}$  assigns a payoff to player  $i$  at each terminal node.
- The *symmetry function*  $T$  is a mapping  $(\cdot)^T$  from choices to choices with the following properties: (1) if  $c \in C_0$ , then  $c^T \in C_0$  and  $p(c) = p(c^T)$ ; (2) if  $c \in C_i$ , then  $c^T \in C_{-i}$ ; (3)  $(c^T)^T = c$  for all  $c$ ; (4) for every information set  $u$  there exists an information set  $u^T$  such that every choice at  $u$  is mapped onto a choice at  $u^T$ ; (5) for every endpoint  $z$  there exists an endpoint  $z^T$  such that if  $z$  is reached by the sequence  $c_1, c_2, \dots, c_k$ , then  $z^T$  is reached by a permutation of  $c_1^T, c_2^T, \dots, c_k^T$ ; and (6)  $r_i(z) = r_{-i}(z^T)$  for every endpoint  $z$ .

As is standard in the literature, I restrict the analysis to games with *perfect recall* (Kuhn [9]). Formally, for each player  $i$ , information sets  $u, v \in U_i$ , choice  $c \in C_u$ , and nodes  $x, y \in v$ , it is assumed that  $c$  comes before  $x$  iff  $c$  comes before  $y$ . A *behavior strategy* of player  $i$  is a mapping that assigns a probability distribution over the set of choices  $C_u$  to every information set  $u \in U_i$ . Let  $B_i$  be the set of all behavior strategies of player  $i$ , and let  $B = B_1 \times B_2$  be the set of all strategy profiles. Given strategy profile  $\mathbf{b} = (b_1, b_2) \in B$ , let  $\mathbb{P}_{\mathbf{b}}(z)$  be the probability that endpoint  $z$  is reached when  $\mathbf{b}$  is played, and let  $R_i(\mathbf{b})$  be the expected payoff to player  $i$  when the players play strategy profile  $\mathbf{b}$ :  $R_i(\mathbf{b}) = \sum_{z \in Z} \mathbb{P}_{\mathbf{b}}(z) \cdot r_i(z)$ . If  $b \in B_1$  is a behavior strategy of player 1 in  $\Gamma$ , then the *symmetric image* of  $b$  is the behavior strategy  $b^T$  of player 2 defined by:  $\beta_u^T(c) := \beta_{u^T}(c^T)$  for each  $u \in U_2$  and  $c \in C_u$ . Observe that the properties of the symmetry function imply that  $R_1(b, b^T) = R_2(b, b^T)$ .

## References

- [1] Bhaskar, Vinit. 1995. On the neutral stability of mixed strategies in asymmetric contests. *Mathematical Social Sciences*, **30**(3), 273–284.
- [2] Bhaskar, Vinit. 1998. Noisy communication and the evolution of cooperation. *Journal of Economic Theory*, **82**(1), 110–131.
- [3] Bolton, Gary E. 1997. The rationality of splitting equally. *Journal of Economic Behavior and Organization*, **32**(3), 365–381.



- 
- [4] Bomze, Immanuel M, & Weibull, Jörgen W. 1995. Does neutral stability imply Lyapunov stability? *Games and Economic Behavior*, **11**(2), 173–192.
- [5] Cressman, Ross. 1997. Local stability of smooth selection dynamics for normal form games. *Mathematical Social Sciences*, **34**(1), 1–19.
- [6] Heller, Y. 2013. Three steps ahead. *mimeo*.
- [7] Kim, Yong-Gwan. 1993. Evolutionary stability in the asymmetric war of attrition. *Journal of theoretical biology*, **161**(1), 13.
- [8] Kim, Yong-Gwan. 1994. Evolutionarily stable strategies in the repeated prisoner’s dilemma. *Mathematical Social Sciences*, **28**(3), 167–197.
- [9] Kuhn, Harold W. 1953. Extensive games and the problem of information. *Contributions to the Theory of Games*, **2**(28), 193–216.
- [10] Leimar, Olof. 1997. Repeated games: a state space approach. *Journal of Theoretical Biology*, **184**(4), 471–498.
- [11] Maynard Smith, J. 1982. *Evolution and the theory of games*. Cambridge University Press.
- [12] Maynard-Smith, J., & Price, G.R. 1973. The Logic of Animal Conflict. *Nature*, **246**, 15.
- [13] Okada, Mr A. 1981. On stability of perfect equilibrium points. *International Journal of Game Theory*, **10**(2), 67–73.
- [14] Samuelson, Larry, & Zhang, Jianbo. 1992. Evolutionary stability in asymmetric games. *Journal of economic theory*, **57**(2), 363–391.
- [15] Sandholm, William H. 2010. Local stability under evolutionary game dynamics. *Theoretical Economics*, **5**(1), 27–50.
- [16] Selten, Reinhard. 1983. Evolutionary stability in extensive two-person games. *Mathematical Social Sciences*, **5**(3), 269–363.
- [17] Taylor, P.D., & Jonker, L.B. 1978. Evolutionary stable strategies and game dynamics. *Mathematical Biosciences*, **40**(1), 145–156.
- [18] van Damme, E. 1987. *Stability and Perfection of Nash Equilibria*. Springer, Berlin.
- [19] Weibull, Jörgen W. 1995. *Evolutionary game theory*. The MIT press.