

# Lecture XV: Axelrod Tournaments

Markus M. Möbius

April 23, 2003

**Readings for this class: Axelrod, chapters 1-3.**

## 1 Introduction

In 1983 Axelrod solicited computer programs from game theorists and psychologists - each program was essentially a strategy of how to play a repeated Prisoner's Dilemma stage game 200 times:

	C	D
C	3,3	0,5
D	5,0	1,1

All 40 entries were matched against each other, and the score of each bilateral match was calculated as the sum of all stage payoffs (i.e.  $\delta = 1$ ). All payoffs were between 200 and 600 since a player can always guarantee at least a payoff

of 1 in a stage game. The average score of each strategy was calculated, and the strategies ranked according to overall success.<sup>1</sup>

Perhaps surprisingly, the winner was the simple TIT FOR TAT rule - a player starts to cooperate and then does whatever the last player did in the previous period.

The reason that TFT did so well was the following:

1. TFT is 'nice' because it never defects first. In fact, all of the best eight strategies were nice. Nice strategies did primarily well because so many of the entries were nice - hence they did particularly well against other nice strategies (getting the maximum payoff of 600).
2. TFT is discriminating enough not to be exploited too badly. For example, 'all C' is easily exploited by 'all D', but TFT can't be as easily duped.
3. TFT is forgiving - bygones are bygones, and one period cooperation can lead to continued cooperation in the future. Many strategies did badly because they were unforgiving. Forgiveness was the main criterion for success amongst the set of 'nice' strategies.

It's ironic that the example on the solicitation letter would have done even better than TFT. This was essentially a TIT for two TATs strategy - a player will only defect if he sees two defections in a row. What's nice about TF2T is that it avoids 'echo effects' where players are caught in an alternate  $(C, D)$  and  $(D, C)$  cycle.

Most submissions were in some way or another variations of TFT. However, rather than making TFT more forgiving (as TF2T) participants devised less nice rules which turned out to do worse on average.

## 2 The Evolutionary Approach

Axelrod uses evolutionary game theory to analyze cooperation in repeated games. This approach is useful because his main interest lies in the question

---

<sup>1</sup>Note, that we wouldn't necessarily expect a strategy to be a NE - after all we know from the folk theorem that there are many SPE implementing a continuum of outcomes. Hence, there is no obvious way to select NE strategies. A winning strategy has to do well on average against other strategies and therefore doesn't have to be Nash (just as in the 2/3 of the average game the winning strategy is typically never part of a Nash equilibrium).

### how does cooperation arise in populations?

Ideally, evolutionary forces should select equilibria with more cooperation over those with less cooperation, and provide some mechanism through which the transition can be accomplished. This mechanism are typically 'mutations' of agents who do better than non-cooperating agents.

In order to analyze evolution in repeated PD games we go through the following steps:

1. Define a suitable evolutionary environment such that our evolutionary concepts (ESS, replicator dynamics etc.) continue to apply in repeated game setting.
2. Show that TFT is evolutionary stable. For simplicity, we use a slightly less demanding concept called collective stability.
3. Realize that lots of strategies are collectively stable, in particular 'All D'.
4. Define the concept of  $q$ -clusters and cluster-invasions and show that TFT is cluster-stable while 'All D' is not. The concept of 'coordinated mutations' is powerful, and has many analogies in real life.

## 3 Evolution and Repeated Games

Our standard evolutionary setting was a mass 1 of agents who are randomly matched all the time and play a stage game against each other. In this case agents will act myopically, even if they have a positive discount factor because the probability of running into the same agent twice is essentially zero (if matching is random and the population is large).

We want to keep the random matching assumption but weaken the requirement that agents' matches break up for sure after each period. Formally, we assume that agents have discount factors  $\delta$  and that there is a probability  $p$  that a match breaks up in every time period (or it breaks up at rate  $p$  if time is continuous). The expected match length is then  $\frac{1}{p}$  - for the Axelrod tournament we would choose  $p = 0.005$  in order to get  $T = 200$  on average. Each agent then discounts payoffs *within the same match* at rate  $\tilde{\delta} = \delta(1 - p)$ . After a match breaks up, agents are instantly rematched. Each period a share  $p$  of agents are rematched.

Previously, agents were labelled by the strategy set (e.g.  $x_i$  was the share of agents using strategy  $s_i$ ). Now, the strategy set is the set of extensive form strategies, i.e.  $s : H \rightarrow \{C, D\}$ .<sup>2</sup>

Now we can define ESS, replicator dynamics etc. exactly as before. To make life easier, we will use a slightly weaker form of ESS, namely collective stability which is essentially NE:

**Definition 1** *A strategy  $s$  is collectively stable if any other strategy  $s'$  (of the infinite extensive form game) does not do better, i.e.:*<sup>3</sup>

$$U(s', s) \leq U(s, s)$$

Collective stability (like NE) captures the idea that no mutation will do better than the original strategy does against itself.<sup>4</sup>

## 4 TFT and Collective Stability

Note, that TFT is NOT a SPE. This can be easily shown using the SPDP. Assume we are in a subgame where both players defected. TFT would lead to continual defection. However, if one player deviates and cooperates, it will lead to an 'echo' of alternating  $(C, D)$  and  $(D, C)$  moves. For sufficiently high discount factor players will do better as a result. Hence TFT is not SPE. However, it is NE and hence collectively stable as the next proposition shows.

Note, that the SPE is hard to translate into an evolutionary setting because after all people are assumed to play fixed rules and hence don't have commitment problems. However, in the repeated PD the NE concept does not allow 'non-credible' threats as NE in the entry game does, for example. The reason is that the only threat players have in PD is in fact playing Nash.

**Proposition 1** *TFT is collectively stable for  $\delta > \frac{2}{3}$ .*

---

<sup>2</sup>There are some potential technical problems because the set of strategies is now infinite (in fact uncountable) whereas before we had to deal only with  $x_C$  and  $x_D$ . Typically, we restrict attention to two strategies - some dominant strategy, and an entrant 'mutant'.

<sup>3</sup>Capital utility is the discounted sum of stage game utilities.

<sup>4</sup>Note, that ESS is stronger because it imposes a second condition.

**Proof:** 'All C' clearly doesn't do better than TFT since TFT is nice. We have to check that 'all D' doesn't do better, and also 'alternate D and C'. Are there any other strategies? NO.

Look at a strategy  $s'$  which contains both C and D along the equilibrium path. Assume there is a string of at least two D's somewhere in that strategy  $s' = (\dots DDXXX\dots)$ . Note, that TFT defects after the first defection. Consider the modified strategy  $s'' = (\dots DXXXX)$ . If  $s''$  does better  $s'$  against TFT then it does not pay to include an additional D in the string of strategy  $s'$ . Hence we can eliminate all double D's from  $s'$  and get a better strategy against TFT. If  $s''$  does not do better than  $s'$  against TFT then we can continue to insert defections to get successively better strategies until we get back to 'all D'. This argument implies that we only have to consider 'all D' or mixed strategies with no consecutive D's.

In the same way we can exclude consecutive C's. This shows that we only have to consider 'all D', and 'alternate D and C'.

$$\begin{aligned} u(\text{TFT}, \text{TFT}) &= \frac{3}{1-\delta} \\ u(\text{all D}, \text{TFT}) &= 5 + \frac{\delta}{1-\delta} \\ u(\text{alternate D and C}, \text{TFT}) &= \frac{5}{1-\delta^2} \end{aligned}$$

We just have to make that none of the two mutations does better than TFT. This will be the case if  $\delta > \max(\frac{1}{2}, \frac{2}{3})$ . QED

## 5 What About Other Strategies?

Any NE is collectively stable. By the folk theorem there is a continuum of SPE (which are obviously NE) and guarantee any IR, feasible payoff (check this area for yourself). In particular:

**Proposition 2** *'All D' is SPE and collectively stable.*

This look like horrible news - how can cooperation ever arise if everything including 'All D' is collectively stable?

## 6 Cluster-stability

For this section assume that  $\delta = .9$  such that TFT is collectively stable. The standard ESS concept of a mutation assumes that mutations are 'small' and hence it's unlikely that a mutant interacts with other mutants.

Axelrod's idea is to consider mutations which are 'small' compared to the population, but whose members nonetheless interact. For example, the senate is pretty large, but two freshmen senators from California might interact sufficiently so that they don't influence the payoffs of other senators very much (just as individual mutations do), but they do affect each other directly a lot.

Formally, we say that agents in a  $q$ -cluster mutation interact with each other with probability  $q$  and interact with everyone else with probability  $1 - q$ .

Recall, that a mutant agent playing TFT gets payoff  $U(\text{TFT}, \text{All D})$  in an 'All D' equilibrium. However, mutants in a  $q$ -cluster get:

$$qU(\text{TFT}, \text{TFT}) + (1 - q)U(\text{TFT}, \text{All D})$$

It can be easily checked that  $q$ -cluster mutants do better than the 'All D' incumbents for  $q > \frac{1}{21} \approx 5$  percent. Hence, cluster mutants don't have to interact a lot to do a lot better than the 'meanies' who are the majority of the population.

**Proposition 3** *TFT is  $q$ -cluster stable for  $q > \frac{1}{21}$  while 'All D' is not  $q$ -cluster stable for any  $q$ .*

### 6.1 Real-World Implications

Hence, a group of meanies is vulnerable to invasion by groups of 'nice' agents such as those playing TFT. In contrast, TFT cannot be invaded neither by individual nor by group mutations. This last theorem is Axelrod's basic insight, and a lot of the examples in the book (senators, World War I trench warfare) revolve around applications of this principle.

**Cooperation can arise quickly and naturally as long as group mutations are possible.**<sup>5</sup>

---

<sup>5</sup>Note, that the intuition is very similar to stochastic stability with local interaction.