

# Stability of defection, optimisation of strategies and the limits of memory in the Prisoner's Dilemma.

Nikoleta E. Glynatsi

Vincent A. Knight

## Abstract

Memory-one strategies are a set of Iterated Prisoner's Dilemma strategies that have been praised for their mathematical tractability and performance against single opponents. This manuscript investigates *best response* memory-one strategies as a multidimensional optimisation problem. Though extortionate memory-one strategies have gained much attention, we demonstrate that best response memory-one strategies do not behave in an extortionate way, and moreover, for memory one strategies to be evolutionarily robust they need to be able to behave in a forgiving way. We also provide evidence that memory-one strategies suffer from their limited memory in multi agent interactions and can be out performed by longer memory strategies.

## 1 Introduction

The Prisoner's Dilemma (PD) is a two player game used in understanding the evolution of co-operative behaviour, formally introduced in [10]. Each player has two options, to cooperate (C) or to defect (D). The decisions are made simultaneously and independently. The normal form representation of the game is given by:

$$S_p = \begin{pmatrix} R & S \\ T & P \end{pmatrix} \quad S_q = \begin{pmatrix} R & T \\ S & P \end{pmatrix} \quad (1)$$

where  $S_p$  represents the utilities of the row player and  $S_q$  the utilities of the column player. The payoffs,  $(R, P, S, T)$ , are constrained by equations (2) and (3). Constraint (2) ensures that defection dominates cooperation and constraint (3) ensures that there is a dilemma; the sum of the utilities for both players is better when both choose to cooperate. The most common values used in the literature are  $(R, P, S, T) = (3, 1, 0, 5)$  [4].

$$T > R > P > S \quad (2)$$

$$2R > T + S \quad (3)$$

The PD is a one shot game, however it is commonly studied in a manner where the history of the interactions matters. The repeated form of the game is called the Iterated Prisoner's Dilemma (IPD) and in the 1980s,

following the work of [5, 6] it attracted the attention of the scientific community. In [5] and [6], the first well known computer tournaments of the IPD were performed. A total of 13 and 63 strategies were submitted respectively in the form of computer code. The contestants competed against each other, a copy of themselves and a random strategy, and the winner was then decided on the average score achieved (not the total number of wins). The contestants were given access to the entire history of a match, however, how many turns of history a strategy would incorporate, refereed to as the *memory size* of a strategy, was a result of the particular strategic decisions made by the author. The winning strategy of both tournaments was the strategy called Tit for Tat and it's success, in both tournaments, came as a surprise. Tit for Tat was a simple, forgiving strategy that opened each interaction by cooperation, but it had managed to defeat far more complicated opponents. Tit for Tat provided evidence that being nice can be advantageous and became the major paradigm for reciprocal altruism.

Another trait of Tit for Tat is that it considers only the previous move of the opponent. These type of strategies are called *reactive* [24] and are a subset of so called *memory-one* strategies, which incorporate both players' latests moves. Reactive and memory-one strategies have been studied thoroughly in, for example [25, 26]. They have gained most of their attention when a certain subset of memory-one strategies was introduced in [28]. In [29] it was stated that "Press and Dyson have fundamentally changed the viewpoint on the Prisoner's Dilemma".

Zero-determinant strategies (ZD) are a special case of memory-one and extortionate strategies. They chose their actions so that a linear relationship is forced between their score and that of the opponent, ensuring that they will always receive at least as much as their opponents. ZD strategies are indeed mathematically unique and are proven to be robust in pairwise interactions. Their true effectiveness in tournament interactions and evolutionary dynamics has been questioned [2, 19, 20].

The purpose of this work is to consider a given memory-one strategy in a similar fashion to [28], however whilst [28] found a way for a player to manipulate a given opponent, this work will consider a multidimensional optimisation approach to identify the best response to a given group of opponents. In particular, this work presents a compact method of identifying the best response memory-one strategy against a given set of opponents.

Further, theoretical and empirical results of this work include:

1. The behaviour of a best response memory-one strategy and whether it behaves extortionately, similar to [28].
2. The factors that make a best response memory-one strategy evolutionary robust.
3. A well designed framework that allows the comparison of an optimal memory one strategy, and a more complex strategy that has a larger memory and was obtained through contemporary reinforcement learning techniques [12].
4. An identification of conditions for which defection is known to be a best response; thus identifying environments where cooperation will not occur.

The source code used in this manuscript has been written in a sustainable manner. It is open source (<https://github.com/Nikoleta-v3/Memory-size-in-the-prisoners-dilemma>) and tested which ensures the validity of the results. It has also been archived and can be found at.

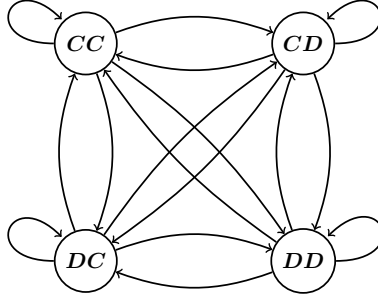


Figure 1: Markov Chain

## 2 The utility

One specific advantage of memory-one strategies is their mathematical tractability. They can be represented completely as an element of  $\mathbb{R}_{[0,1]}^4$ . This originates from [24] where it is stated that if a strategy is concerned with only the outcome of a single turn then there are four possible ‘states’ the strategy could be in;  $CC, CD, DC, CC$ . Therefore, a memory-one strategy can be denoted by the probability vector of cooperating after each of these states;  $p = (p_1, p_2, p_3, p_4) \in \mathbb{R}_{[0,1]}^4$ . In an IPD match two memory-one strategies are moving from state to state at each turn with a given probability. This exact behaviour can be modeled as a stochastic process, and more specifically as a Markov chain (Figure 1). The corresponding transition matrix  $M$  of Figure 1 is given in (4),

$$M = \begin{bmatrix} p_1 q_1 & p_1 (-q_1 + 1) & q_1 (-p_1 + 1) & (-p_1 + 1) (-q_1 + 1) \\ p_2 q_3 & p_2 (-q_3 + 1) & q_3 (-p_2 + 1) & (-p_2 + 1) (-q_3 + 1) \\ p_3 q_2 & p_3 (-q_2 + 1) & q_2 (-p_3 + 1) & (-p_3 + 1) (-q_2 + 1) \\ p_4 q_4 & p_4 (-q_4 + 1) & q_4 (-p_4 + 1) & (-p_4 + 1) (-q_4 + 1) \end{bmatrix} \quad (4)$$

The long run steady state probability vector  $v$  is the solution to  $vM = v$ . The stationary vector  $v$  can be combined with the payoff matrices of (1) and the expected payoffs for each player can be estimated without simulating the actual interactions. More specifically, the utility for a memory-one strategy  $p$  against an opponent  $q$ , denoted as  $u_q(p)$ , is defined by,

$$u_q(p) = v \cdot (R, S, T, P). \quad (5)$$

The first theoretical result of this manuscript is presented in Theorem 1. Theorem 1 states that  $u_q(p)$  is given by a ratio of two quadratic forms [17]. To the authors knowledge our work is the first to explore the form of  $u_q(p)$ .

**Theorem 1.** *The expected utility of a memory-one strategy  $p \in \mathbb{R}_{[0,1]}^4$  against a memory-one opponent  $q \in \mathbb{R}_{[0,1]}^4$ , denoted as  $u_q(p)$ , can be written as a ratio of two quadratic forms:*

$$u_q(p) = \frac{\frac{1}{2}pQp^T + cp + a}{\frac{1}{2}p\bar{Q}p^T + \bar{c}p + \bar{a}}, \quad (6)$$

where  $Q, \bar{Q} \in \mathbb{R}^{4 \times 4}$  are square matrices whose diagonal elements are all equal to zero, and are defined by the transition probabilities of the opponent  $q_1, q_2, q_3, q_4$  as follows:

$$Q = \begin{bmatrix} 0 & -(q_1 - q_3)(q_2 - 5q_4 - 1) & q_3(q_1 - q_2) & -5q_3(q_1 - q_4) \\ -(q_1 - q_3)(q_2 - 5q_4 - 1) & 0 & (q_2 - q_3)(q_1 - 3q_4 - 1) & (q_3 - q_4)(5q_1 - 3q_2 - 2) \\ q_3(q_1 - q_2) & (q_2 - q_3)(q_1 - 3q_4 - 1) & 0 & 3q_3(q_2 - q_4) \\ -5q_3(q_1 - q_4) & (q_3 - q_4)(5q_1 - 3q_2 - 2) & 3q_3(q_2 - q_4) & 0 \end{bmatrix}, \quad (7)$$

$$\bar{Q} = \begin{bmatrix} 0 & -(q_1 - q_3)(q_2 - q_4 - 1) & (q_1 - q_2)(q_3 - q_4) & (q_1 - q_4)(q_2 - q_3 - 1) \\ -(q_1 - q_3)(q_2 - q_4 - 1) & 0 & (q_2 - q_3)(q_1 - q_4 - 1) & (q_1 - q_2)(q_3 - q_4) \\ (q_1 - q_2)(q_3 - q_4) & (q_2 - q_3)(q_1 - q_4 - 1) & 0 & -(q_2 - q_4)(q_1 - q_3 - 1) \\ (q_1 - q_4)(q_2 - q_3 - 1) & (q_1 - q_2)(q_3 - q_4) & -(q_2 - q_4)(q_1 - q_3 - 1) & 0 \end{bmatrix}. \quad (8)$$

$c$  and  $\bar{c} \in \mathbb{R}^{4 \times 1}$  are similarly defined by:

$$c = \begin{bmatrix} q_1(q_2 - 5q_4 - 1) \\ -(q_3 - 1)(q_2 - 5q_4 - 1) \\ -q_1q_2 + q_2q_3 + 3q_2q_4 + q_2 - q_3 \\ 5q_1q_4 - 3q_2q_4 - 5q_3q_4 + 5q_3 - 2q_4 \end{bmatrix}, \quad (9)$$

$$\bar{c} = \begin{bmatrix} q_1(q_2 - q_4 - 1) \\ -(q_3 - 1)(q_2 - q_4 - 1) \\ -q_1q_2 + q_2q_3 + q_2 - q_3 + q_4 \\ q_1q_4 - q_2 - q_3q_4 + q_3 - q_4 + 1 \end{bmatrix}. \quad (10)$$

and  $a = -q_2 + 5q_4 + 1$  and  $\bar{a} = -q_2 + q_4 + 1$ .

The proof of Theorem 1 is given in Appendix A.

Numerical simulations have been carried out to validate the formulation of  $u_q(p)$  as a quadratic ratio. The simulated utility, which is denoted as  $U_q(p)$ , has been calculated using [1] an open source research framework for the study of the IPD <sup>1</sup>. For smoothing the simulated results the simulated utility has been estimated in a tournament of 500 turns and 200 repetitions. Figure 2 shows that the formulation of Theorem 1 successfully captures the simulated behaviour.

Theorem 1 can be extended to consider multiple opponents. The IPD is commonly studied in tournaments and/or Moran Processes where a strategy interacts with a number of opponents. The payoff of a player in such interactions is given by the average payoff the player received against each opponent. More specifically the expected utility of a memory-one strategy against a  $N$  number of opponents is given by Theorem 2.

**Theorem 2.** *The expected utility of a memory-one strategy  $p \in \mathbb{R}_{[0,1]}^4$  against a group of opponents  $q^{(1)}, q^{(2)}, \dots, q^{(N)}$ , denoted as  $\frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p)$ , is given by:*

---

<sup>1</sup>The project is described in [18].

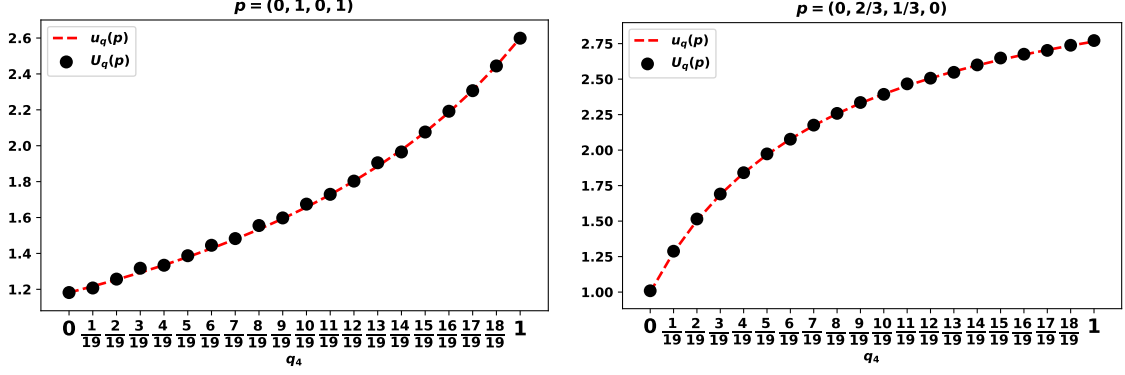


Figure 2: Simulated and analytically calculated utility for  $p = (0, 1, 0, 1)$  and  $p = (0, \frac{2}{3}, \frac{1}{3}, 0)$  against  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, q_4)$  for  $q_4 \in \{0, \frac{1}{19}, \frac{2}{19}, \dots, \frac{18}{19}, 1\}$ .

$$\frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p) = \frac{1}{N} \frac{\sum_{i=1}^N (\frac{1}{2} p Q^{(i)} p^T + c^{(i)} p + a^{(i)}) \prod_{\substack{j=1 \\ j \neq i}}^N (\frac{1}{2} p \bar{Q}^{(j)} p^T + \bar{c}^{(j)} p + \bar{a}^{(j)})}{\prod_{i=1}^N (\frac{1}{2} p \bar{Q}^{(i)} p^T + \bar{c}^{(i)} p + \bar{a}^{(i)})}. \quad (11)$$

The proof of Theorem 2 is a straightforward algebraic manipulation.

Similar to the previous result, the formulation of Theorem 2 is validated using numerical simulations. The simulated and formulated utilities of strategies in a tournament of 10 opponents as described in [29] are a match, Figure 3.

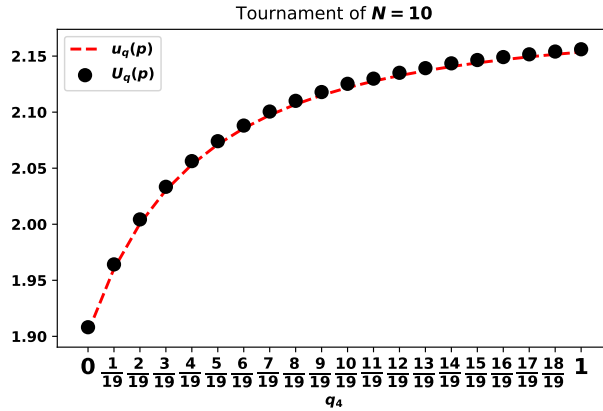


Figure 3: The utilities of memory-one strategies  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, p_4)$  for  $p_4 \in \{0, \frac{1}{19}, \frac{2}{19}, \dots, \frac{18}{19}, 1\}$  against the 10 memory-one strategies used in [29].

Furthermore, using the list of strategies from [29] we check whether the utility against a group of strategies could be captured by the utility against the mean opponent. Thus whether condition 12 holds. However, condition (12) fails as shown in Figure 4.

$$\frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p) = u_{\frac{1}{N} \sum_{i=1}^N q^{(i)}}(p), \quad (12)$$

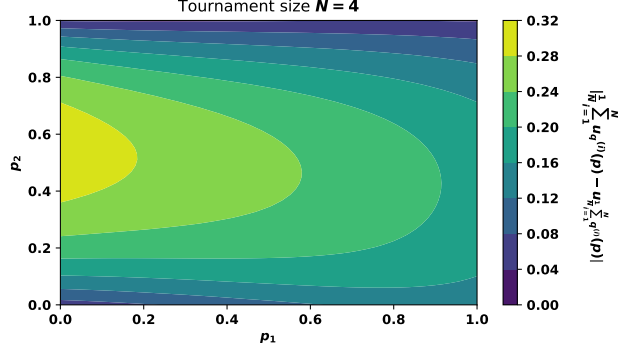


Figure 4: The difference between the average utility and against the utility against the average player of the strategies in [29]. A positive difference indicates that the condition (12) does not hold.

Two theoretical results have been presented so far. The formulation of Theorem 2 which allows for the utility of a memory-one strategy against any number of opponents to be estimated without simulating the interactions is the main results used in this manuscript. In Section 3 it is used to define best response memory-one strategies and explore the conditions under which defection dominates cooperation.

### 3 Best responses to memory-one players

This section focused on best responses and more specifically *best response memory-one* strategies. A *best response* is the strategy which corresponds to the most favorable outcome [31], thus a best response memory-one corresponds to a strategy  $p^*$  for which (11) is maximised. This is considered as a multi dimensional optimisation problem given by,

$$\begin{aligned} \max_p : & \sum_{i=1}^N u_q^{(i)}(p) \\ \text{such that : } & p \in \mathbb{R}_{[0,1]} \end{aligned} \quad (13)$$

The decision variable is the vector  $p$ , the solitary constraint is that  $p \in \mathbb{R}_{[0,1]}^4$  and the objective function is (11).

Optimising this particular ratio of quadratic forms is not trivial. It can be verified empirically for the case of a single opponent that there exist at least one point for which the definition of concavity does not hold. Some results are known for non concave ratios of quadratic forms [7, 9], however, in these works it's assumed that either both the numerator and the denominator of the fractional problem are concave or that the denominator is greater than zero. Both assumptions fail here as stated in Theorem 3.

**Theorem 3.** *The utility of a player  $p$  against an opponent  $q$ ,  $u_q(p)$  given by (6), is not concave. Furthermore neither the numerator or the denominator of (6), are concave.*

Proof is given in Appendix A.2.

The non concavity of  $u(p)$  indicates multiple local optimal points. The approach taken here is to introduce a compact way of constructing the candidate set of all local optimal points. Once the set is defined the point that maximises (11) corresponds to the best response strategy. The problem considered is bounded because  $p \in \mathbb{R}_{[0,1]}^4$ . Thus, the candidate solutions will exist either at the boundaries of the feasible solution space, or within that space. The method of Lagrange Multipliers [8] and Karush-Kuhn-Tucker conditions [11] are based on this.

This approach allow us to define the best response memory-one strategy to a group of opponents, Lemma 4.

**Lemma 4.** *The optimal behaviour of a memory-one strategy player  $p^* \in \mathbb{R}_{[0,1]}^4$  against a set of  $N$  opponents  $\{q^{(1)}, q^{(2)}, \dots, q^{(N)}\}$  for  $q^{(i)} \in \mathbb{R}_{[0,1]}^4$  is established by:*

$$p^* = \operatorname{argmax} \left( \sum_{i=1}^N u_q(p) \right), \quad p \in S_q.$$

The set  $S_q$  is defined as all the possible combinations of:

$$S_q = \left\{ p \in \mathbb{R}^4 \left| \begin{array}{l} \bullet \quad p_j \in \{0, 1\} \quad \text{and} \quad \frac{d}{dp_k} \sum_{i=1}^N u_q^{(i)}(p) = 0 \quad \text{for all } j \in J \quad \& \quad k \in K \quad \text{for all } J, K \\ \bullet \quad p \in \{0, 1\}^4 \end{array} \right. \right\}.$$

*where  $J \cap K = \emptyset$  and  $J \cup K = \{1, 2, 3, 4\}$ .*

The proof is given in the Appendix A.

Note that there is no immediate way to find the zeros of  $\frac{d}{dp} \sum_{i=1}^N u_q(p)$ ;

$$\begin{aligned} \frac{d}{dp} \sum_{i=1}^N u_q^{(i)}(p) &= \\ &= \sum_{i=1}^N \frac{\left( pQ^{(i)} + c^{(i)} \right) \left( \frac{1}{2}p\bar{Q}^{(i)}p^T + \bar{c}^{(i)}p + \bar{a}^{(i)} \right) - \left( p\bar{Q}^{(i)} + \bar{c}^{(i)} \right) \left( \frac{1}{2}pQ^{(i)}p^T + c^{(i)}p + a^{(i)} \right)}{\left( \frac{1}{2}p\bar{Q}^{(i)}p^T + \bar{c}^{(i)}p + \bar{a}^{(i)} \right)^2} \end{aligned} \quad (14)$$

For  $\frac{d}{dp} \sum_{i=1}^N u_q(p)$  to equal zero then:

$$\sum_{i=1}^N \left( \left( pQ^{(i)} + c^{(i)} \right) \left( \frac{1}{2}p\bar{Q}^{(i)}p^T + \bar{c}^{(i)}p + \bar{a}^{(i)} \right) - \left( p\bar{Q}^{(i)} + \bar{c}^{(i)} \right) \left( \frac{1}{2}pQ^{(i)}p^T + c^{(i)}p + a^{(i)} \right) \right) = 0, \quad \text{while} \quad (15)$$

$$\sum_{i=1}^N \frac{1}{2}p\bar{Q}^{(i)}p^T + \bar{c}^{(i)}p + \bar{a}^{(i)} \neq 0. \quad (16)$$

Constructing the subset  $S_q$  is analytically possible. The points for any or all of  $p_i \in \{0, 1\}$  for  $i \in \{1, 2, 3, 4\}$

are trivial. Finding the roots of the partial derivatives  $\frac{d}{dp} \sum_{i=1}^N u_q(p)$  which are a set of polynomials of equations (15) is feasible using resultant theory. Resultant theory [16] allow us to solve systems of polynomials by the calculation of a resultant.

However, for large systems these quickly become intractable and numerical methods taking advantage of the structure will be used which are described in Section 4. The rest of the section focuses on an immediate theoretical result from Lemma 4.

### 3.1 Stability of defection

An immediate result from Lemma 4 can be obtained by evaluating the sign of the derivative (14) at  $p = (0, 0, 0, 0)$ . If at that point the derivative is negative, then the utility of the player is maximum at that point and it will only decrease if the player were to change their behaviour. Thus, defection is the best response.

**Lemma 5.** *In a tournament of  $N$  players where  $q^{(i)} = (q_1^{(i)}, q_2^{(i)}, q_3^{(i)}, q_4^{(i)})$  defection is a best response if the transition probabilities of the opponents satisfy the condition (18).*

$$\sum_{i=1}^N (c^{(i)T} \bar{a}^{(i)} - \bar{c}^{(i)T} a^{(i)}) \leq 0 \quad (17)$$

while,

$$\sum_{i=1}^N \bar{a}^{(i)} \neq 0 \quad (18)$$

*Proof.* For defection to be a best response the derivative of the utility at the point  $p = (0, 0, 0, 0)$  must be negative. This would indicate that the utility function is only declining from that point onwards.

Substituting  $p = (0, 0, 0, 0)$  in equation (14) gives:

$$\sum_{i=1}^N \frac{(c^{(i)T} \bar{a}^{(i)} - \bar{c}^{(i)T} a^{(i)})}{(\bar{a}^{(i)})^2} \quad (19)$$

The sign of the numerator  $\sum_{i=1}^N (c^{(i)T} \bar{a}^{(i)} - \bar{c}^{(i)T} a^{(i)})$  can vary based on the transition probabilities of the opponents. The denominator can not be negative, and otherwise is always positive. Thus the sign of the derivative is negative if and only if  $\sum_{i=1}^N (c^{(i)T} \bar{a}^{(i)} - \bar{c}^{(i)T} a^{(i)}) \leq 0$ .  $\square$

In an environment where defection is the best response the average payoff of a defector is always higher than any other strategy can achieve. If we consider a setting where in each the prevalence of each type of strategy was determined by that strategy's success in the previous round, then in a population such that (18) holds, defection would prevail; thus cooperation would never occur.



This is demonstrated in Figures 5 and 6.

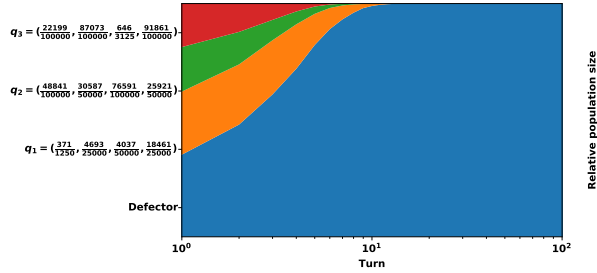


Figure 5: For given opponents condition (18) and Defector takes over the population.

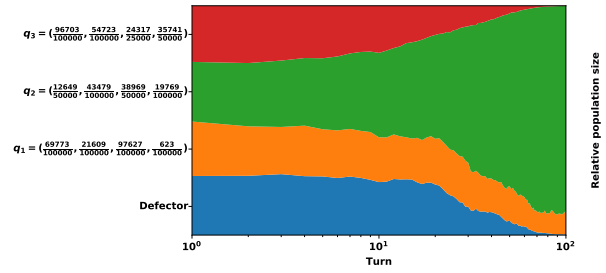


Figure 6: For given opponents condition (18) fails and Defector does not take over the population.

## 4 Numerical experiments

This section focuses on a series of numerical experiments in order to gain a better understanding of memory-one strategies, their behaviour, robustness and limitations. The problems that will be described in this section are solved heuristically using Bayesian optimisation [23].

Bayesian optimisation is a global optimisation algorithm that has proven to outperform many other popular algorithms [15]. The algorithm builds a bayesian understanding of the objective function and it is well suited to the multiple local optimas in the described search area of this work. Differential evolution [30] was also considered, however it was not selected due to Bayesian being computationally more efficient.

For example consider the problem of (13) where  $N = 2$ . Figure 7 illustrates the change of the utility function over iterations of the algorithm. The default number of iterations that has been used in this work is 60. After 60 calls the convergence of the utility is checked. If the optimised utility has changed in the last 10% iterations then a further 20 iterations are considered.

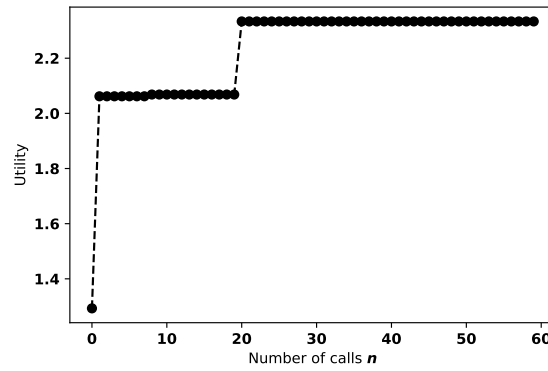


Figure 7: Utility over time of calls using Bayesian optimisation. The opponents are  $q^{(1)} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})$  and  $q^{(2)} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ . The best response obtained is  $p^* = (0.0, \frac{11}{50}, 0.0, 0.0)$

The rest of this section organized as follows. In Section 4.1 is used to estimate a large number of best response memory one strategies and explore whether they behave in an extortionate way. In Section 4.2 similar, a

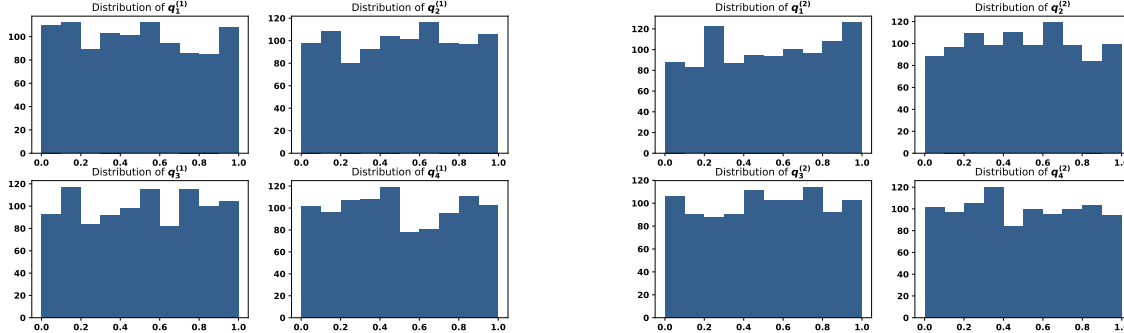
large set of best responses are estimated by this time self interactions are included. Finally in Section 4.3, for a large number of opponents we compare the utility of best response memory one and longer memory strategies.

#### 4.1 Best response memory-one strategies for $N = 2$

As briefly discussed in Section 1 zero-determinant strategies have been praised for their robustness against a single opponent. By forcing a linear relationship between the scores zero-determinant strategies can always receive a higher, or in the case of mutual defection, the same payoff as their opponents. In IPD tournaments the winner is decided on the average score a strategy received and not by wins. Thus winning against an opponent does not guarantee a strategy's success.

We argue that by trying to exploit their opponents zero-determinant strategies suffer in multi opponent interaction where the payoffs matter. In comparison, best response memory-one strategies utilise their behaviour to gain the most from their interactions. The aim of this section is to understand whether best responses behave in an extortionate way, similarly to zero determinants. To estimate a strategy's extortionate behaviour the SSE method as described in [Knight 2019] is used. SSE is defined as the error, how far, a strategy is from behaving extortionate. For example, a high SSE implies that a strategy is not extortionate.

A large data set of best response memory-one strategies when  $N = 2$  has been generated and is available here. trials corresponding to 10000 different best responses. For each trial a set of 2 opponents is randomly generated, the memory-one best response against them is estimated. Though the probabilities  $q_i$  of the opponents are randomly generated, Figures 8a and 8b, show that they are uniformly distributed over the trial. Thus, the full space of possible opponents has been covered.



(a) Distributions of first opponents' probabilities.

(b) Distributions of second opponents' probabilities.

The SSE method has been applied to the data set and it's distribution is shown in Figure 11 alongside a statistics summary in Table 2.

The distribution of SSE is skewed to left indicating that the best response does exhibit extortionate behaviour, however, the best response is not uniformly extortionate. A positive measure of skewness and kurtosis indicate a heavy tail to the right, therefore, in several cases the strategy is not trying to extort the opponent. To conclude, a best response strategy utilises its performance by behaving in a more adaptable way than zero-determinant strategies. This analysis is extended to an evolutionary setting.

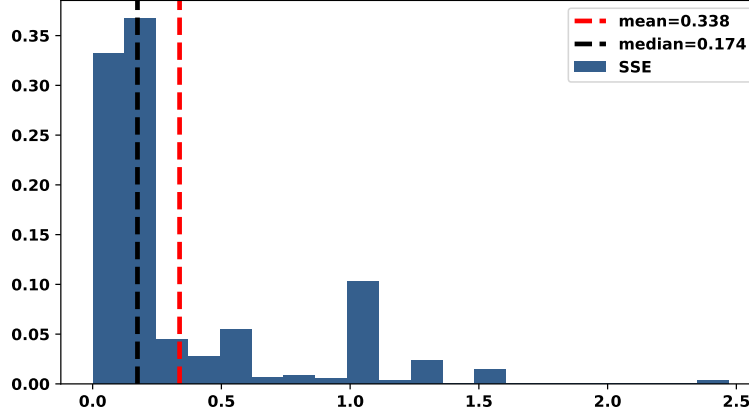


Figure 9: Distribution of SSE for memory-one best responses, when  $N = 2$ .

SSE	
count	1000.00000
mean	0.33762
std	0.39667
min	0.00000
5%	0.02078
25%	0.07597
50%	0.17407
95%	1.05943
max	2.47059
median	0.17407
skew	1.87231
kurt	3.60029

Table 1: Summary statistics SSE of best response memory one strategies included tournaments of  $N = 2$ .

## 4.2 Memory-one best responses in evolutionary dynamics

The IPD is commonly studied in evolutionary processes where the strategies that compose the population can adapt and change their behaviour based on the outcomes of their interactions at each generation. In these processes self interactions are key.

Self interactions can be incorporated in the formulation that has been used in this paper. The utility of a memory-one strategy in an evolutionary setting is given by,

$$\frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p) + u_p(p). \quad (20)$$

and respectively the optimisation problem of (13) is now re written as,

$$\begin{aligned} \max_p : & \frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p) + u_p(p) \\ \text{such that : } & p \in \mathbb{R}_{[0,1]} \end{aligned} \quad (21)$$

Solving this can done using *best response dynamics* as detailed in Algorithm 1. Best response dynamics are commonly used in evolutionary game theory. They represent a class of strategy updating rules, where players

strategies in the next round are determined by their best responses to some subset of the population.

```

 $p^{(t)} \leftarrow (1, 1, 1, 1);$ 
while  $p^{(t)} \neq p^{(t-1)}$  do
     $p^{(t+1)} = \operatorname{argmax} \frac{1}{N} \sum_{i=1}^N u_q^{(i)}(p^{(t+1)}) + u_p^{(t)}(p^{(t+1)});$ 
end

```

**Algorithm 1:** Best response dynamics Algorithm

Algorithm 1 starts by setting an initial solution  $p^{(1)} = (1, 1, 1, 1)$  and repeatedly finds a strategy that maximises (21) using the numerical approach described in Section 4. The algorithm stops when cycle (a sequence of iterated evaluated points) is detected. A numerical example is given in Figure 10.

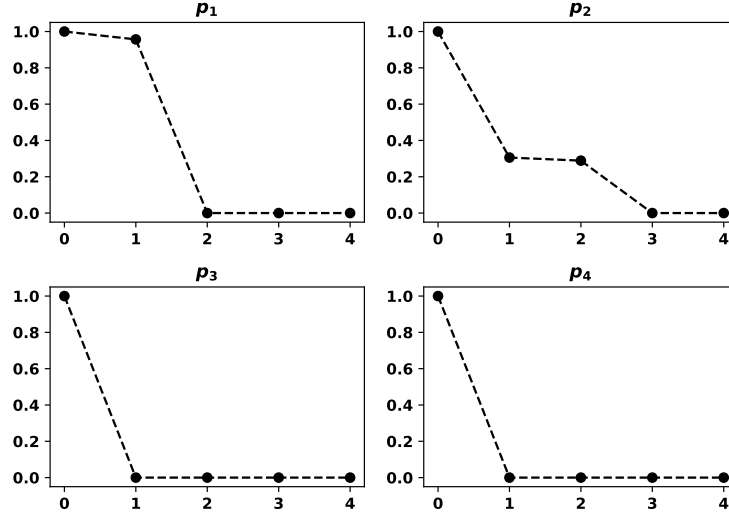


Figure 10: Best response dynamics with  $N = 2$ . More specifically, for  $q^{(1)} = (0.2360, 0.1031, 0.3960, 0.1549)$  and  $q^{(2)} = (0.0665, 0.4015, 0.9179, 0.8004)$ .

For each of the 1000 pairs of opponents, from Section 4.1, the best response in an evolutionary setting has also been obtained. The distribution of SSE for the best response using evolutionary dynamics is given in Figure 11 and a statistical summary of it's distribution in Table 2.

Similarly to the results of Section 4.1, the evolutionary best response strategy does not behave uniformly extortionate. A larger value of both the kurtosis and the skewness of the SSE distribution indicates that an evolutionary best response is more adaptable than the equivalent best response. The difference between the strategies is further explored. Figure 12 compares the tournament and evolutionary best responses. Table 3 details that no statistically significant differences have been found. However, from Figure 12, it seems that evolutionary best response is less likely to forgive after mutual defection whilst it is more likely to forgive after being tricked (higher  $p_2$  median).

### 4.3 Longer memory best response

The effectiveness of memory have been studied thoroughly in literature and evidence where longer memories performed better and were more robustness. This section focuses on proving that short memory strategies

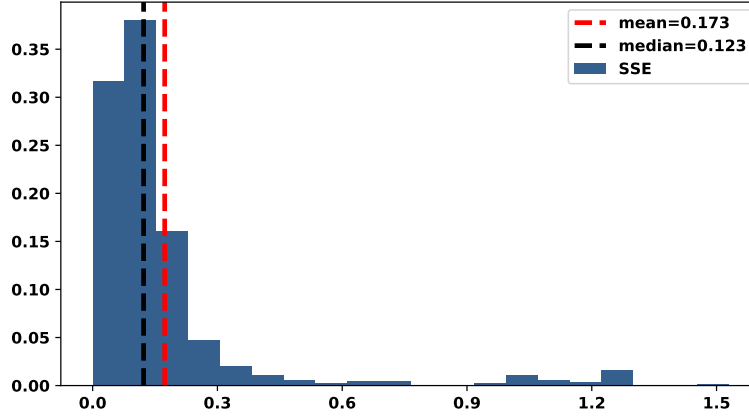


Figure 11: Distribution of SSE of best response memory-one strategies in evolutionary settings, when  $N = 2$ .

SSE	
count	1000.00000
mean	0.17326
std	0.23489
min	0.00001
5%	0.01497
25%	0.05882
50%	0.12253
95%	0.67429
max	1.52941
median	0.12253
skew	3.41839
kurt	11.92339

Table 2: Summary statistics SSE of best response memory-one strategies in evolutionary settings, when  $N = 2$ .

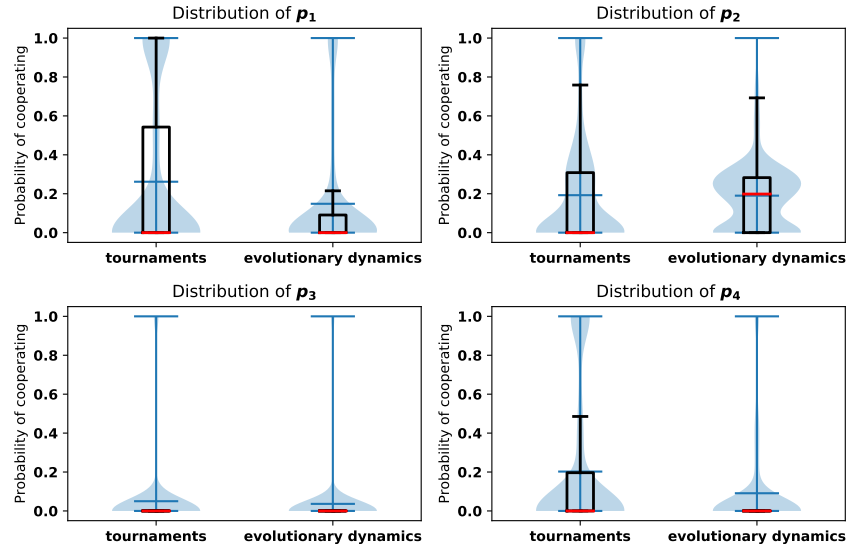


Figure 12: Distributions of  $p^*$  for both best response and evo memory-one strategies.

	Best Response Median in:	ournament	Evolutionary Settings	p-values
0	Distribution $p_1$	0.0	0.00000	0.0
1	Distribution $p_2$	0.0	0.19847	0.0
2	Distribution $p_3$	0.0	0.00000	0.0
3	Distribution $p_4$	0.0	0.00000	0.0

Table 3: A non parametric test, Wilcoxon Rank Sum, has been performed to tests the difference in the medians. A non parametric test is used because is evident that the data are skewed.

have limitations. More specifically, we present several empirical results that show that longer memories strategies can indeed perform better in cases of  $N = 2$ .

Similarly to Sections, a random number of 2 opponents is selcted and then the best memeory one strategy is calculated, algosidned we calucated train to be etter.

This is achieved by comparing the performance of an optimised memory-one strategy to that of a trained long memory-one.

The longer memory strategy selected is a strategy called *Gambler*, introduced and discussed in [12]. A Gambler strategy makes probabilistic decisions based on the opponent's  $n_1$  first moves, the opponent's  $m_1$  last moves and the player's  $m_2$  last moves. This manuscript considers  $n_1 = 2, m_1 = 1$  and  $m_2 = 1$ . By considering the opponent's first two moves, the opponents last move and the player's last move, there are only 16 possible outcomes that can occur. Gambler also makes a probabilistic decision of cooperating in the first move. Combining these Gambler is a function  $f : \{C, D\}^{4 \cup 1} \rightarrow (0, 1)_{\mathbb{R}}$

So this can be hard coded as an element of  $[0, 1]_{\mathbb{R}}^{16+1}$  one probability for each outcome plus the opening move. Thus in comparison to (), finding and optimal Gambler is a 17 dimensional problem, which is solved numerically using Bayesian optimisation. The optimisation problem is given by:

$$\begin{aligned} \max_p : \sum_{i=1}^N U_q^{(i)}(F) \\ \text{such that : } F \in \mathbb{R}_{[0,1]}^{17} \end{aligned} \quad (22)$$

A graphical representation of Gambler and more specifically Gambler( $n_1 = 2, m_1 = 1, m_2 = 1$ ) is given by Figure 13.

Bayesian optimisation is used to numerically solve (22). Similarly to the other experiments, two random opponents are generated and the trained Gambler as well as the best response memory-one are recorded for each trial. A total of 89 trials have been recorded. The utility of both strategies for each trial is estimated, Figure 14.

Though Gambler has an infinite memory (in order to remember the opening moves of the opponent) the information the strategy considers is not significantly larger than memory-one strategies. Even so, it is evident from Figure 14 that Gambler will always performs the same or better than a best response memory one strategy, thus having a longer memory is beneficial.

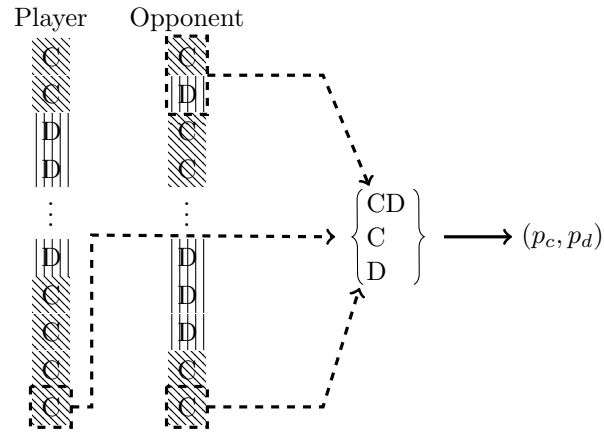


Figure 13: Graphical representation of Gambler.

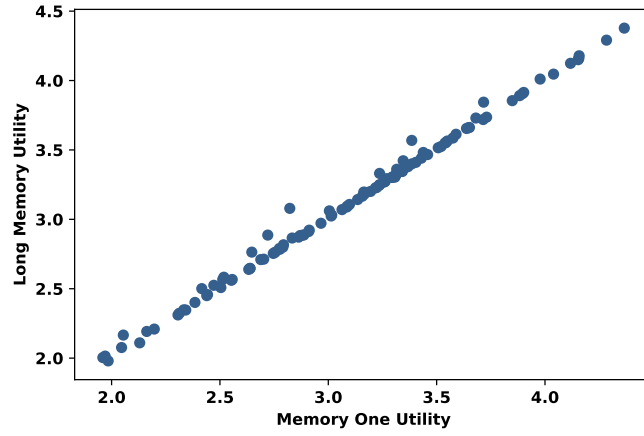


Figure 14: Utilities of Gambler and best response memory-one strategies for 89 different pair of opponents.

## 5 Conclusion

## 6 Acknowledgements

A variety of software libraries have been used in this work:

- The Axelrod library for IPD simulations [1].
- The Scikit-optimize library for an implementation of Bayesian optimisation [13].
- The Matplotlib library for visualisation [14].
- The SymPy library for symbolic mathematics [22].
- The Numpy library for data manipulation [32].

## References

- [1] The Axelrod project developers . Axelrod: 4.4.0, April 2016.
- [2] Christoph Adami and Arend Hintze. Evolutionary instability of zero-determinant strategies demonstrates that winning is not everything. *Nature communications*, 4:2193, 2013.
- [3] Howard Anton and Chris Rorres. *Elementary Linear Algebra: Applications Version*. Wiley, eleventh edition, 2014.
- [4] R Axelrod and WD Hamilton. The evolution of cooperation. *Science*, 211(4489):1390–1396, 1981.
- [5] Robert Axelrod. Effective choice in the prisoner’s dilemma. *The Journal of Conflict Resolution*, 24(1):3–25, 1980.
- [6] Robert Axelrod. More effective choice in the prisoner’s dilemma. *The Journal of Conflict Resolution*, 24(3):379–403, 1980.
- [7] Amir Beck and Marc Teboulle. A convex optimization approach for minimizing the ratio of indefinite quadratic functions over an ellipsoid. *Mathematical Programming*, 118(1):13–35, 2009.
- [8] Dimitri P Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic press, 2014.
- [9] Hongyan Cai, Yanfei Wang, and Tao Yi. An approach for minimizing a quadratically constrained fractional quadratic problem with application to the communications over wireless channels. *Optimization Methods and Software*, 29(2):310–320, 2014.
- [10] Merrill M. Flood. Some experimental games. *Management Science*, 5(1):5–26, 1958.
- [11] Giorgio Giorgi, Bienvenido Jiménez, and Vicente Novo. Approximate karush—kuhn—tucker condition in multiobjective optimization. *J. Optim. Theory Appl.*, 171(1):70–89, October 2016.
- [12] Marc Harper, Vincent Knight, Martin Jones, Georgios Koutsououlos, Nikoleta E. Glynatsi, and Owen Campbell. Reinforcement learning produces dominant strategies for the iterated prisoners dilemma. *PLOS ONE*, 12(12):1–33, 12 2017.



- [13] Tim Head, MechCoder, Gilles Louppe, Iaroslav Shcherbatyi, fcharras, Z Vincius, cmmalone, Christopher Schrder, nel215, Nuno Campos, Todd Young, Stefano Cereda, Thomas Fan, rene rex, Kejia (KJ) Shi, Justus Schwabedal, carlosdanielcsantos, Hvass-Labs, Mikhail Pak, SoManyUsernamesTaken, Fred Callaway, Loc Estve, Lilian Besson, Mehdi Cherti, Karlson Pfannschmidt, Fabian Linzberger, Christophe Cauet, Anna Gut, Andreas Mueller, and Alexander Fabisch. `scikit-optimize/scikit-optimize: v0.5.2`, March 2018.
- [14] J. D. Hunter. Matplotlib: A 2D graphics environment. *Computing In Science & Engineering*, 9(3):90–95, 2007.
- [15] Donald R Jones. A taxonomy of global optimization methods based on response surfaces. *Journal of global optimization*, 21(4):345–383, 2001.
- [16] Gubjorn Jonsson and Stephen Vavasis. Accurate solution of polynomial equations using macaulay resultant matrices. *Mathematics of computation*, 74(249):221–262, 2005.
- [17] Jeremy Kepner and John Gilbert. *Graph algorithms in the language of linear algebra*. SIAM, 2011.
- [18] Vincent Knight, Owen Campbell, Marc Harper, Karol Langner, James Campbell, Thomas Campbell, Alex Carney, Martin Chorley, Cameron Davidson-Pilon, Kristian Glass, Tomáš Ehrlich, Martin Jones, Georgios Koutsououlos, Holly Tibble, Müller Jochen, Geraint Palmer, Paul Slavin, Timothy Standen, Luis Visintini, and Karl Molden. An open reproducible framework for the study of the iterated prisoner’s dilemma. 1(1), 2016.
- [19] Vincent Knight, Marc Harper, Nikoleta E. Glynatsi, and Owen Campbell. Evolution reinforces cooperation with the emergence of self-recognition mechanisms: An empirical study of strategies in the moran process for the iterated prisoners dilemma. *PLOS ONE*, 13(10):1–33, 10 2018.
- [20] Christopher Lee, Marc Harper, and Dashiell Fryer. The art of war: Beyond memory-one strategies in population games. *PLOS ONE*, 10(3):1–16, 03 2015.
- [21] Jiawei Li, Philip Hingston, and Graham Kendall. Engineering design of strategies for winning iterated prisoner’s dilemma competitions. *IEEE Transactions on Computational Intelligence and AI in Games*, 3(4):348–360, 2011.
- [22] A. Meurer, C. P. Smith, M. Paprocki, O. Čertík, S. B. Kirpichev, M. Rocklin, A. Kumar, S. Ivanov, J. K. Moore, S. Singh, T. Rathnayake, S. Vig, B. E. Granger, R. P. Muller, F. Bonazzi, H. Gupta, S. Vats, F. Johansson, F. Pedregosa, M. J. Curry, A. R. Terrel, Š. Roučka, A. Saboo, I. Fernando, S. Kulal, R. Cimrman, and A. Scopatz. Sympy: symbolic computing in python. *PeerJ Computer Science*, 3, 2017.
- [23] J. Moćkus. On bayesian methods for seeking the extremum. In G. I. Marchuk, editor, *Optimization Techniques IFIP Technical Conference Novosibirsk, July 1–7, 1974*, pages 400–404, Berlin, Heidelberg, 1975. Springer Berlin Heidelberg.
- [24] Martin Nowak and Karl Sigmund. Game-dynamical aspects of the prisoner’s dilemma. *Applied Mathematics and Computation*, 30(3):191–213, 1989.
- [25] Martin Nowak and Karl Sigmund. The evolution of stochastic strategies in the prisoner’s dilemma. *Acta Applicandae Mathematicae*, 20(3):247–265, 1990.
- [26] Martin Nowak and Karl Sigmund. A strategy of win-stay, lose-shift that outperforms tit-for-tat in the prisoner’s dilemma game. *Nature*, 364(6432):56, 1993.
- [27] Hisashi Ohtsuki, Christoph Hauert, Erez Lieberman, and Martin A Nowak. A simple rule for the evolution of cooperation on graphs and social networks. *Nature*, 441(7092):502, 2006.

- [28] William H. Press and Freeman J. Dyson. Iterated prisoners dilemma contains strategies that dominate any evolutionary opponent. *Proceedings of the National Academy of Sciences*, 109(26):10409–10413, 2012.
- [29] Alexander J. Stewart and Joshua B. Plotkin. Extortion and cooperation in the prisoners dilemma. *Proceedings of the National Academy of Sciences*, 109(26):10134–10135, 2012.
- [30] Rainer Storn and Kenneth Price. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization*, 11(4):341–359, 1997.
- [31] Steve Tadelis. *Game theory: an introduction*. Princeton University Press, 2013.
- [32] S. Walt, S. C. Colbert, and G. Varoquaux. The NumPy array: a structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22–30, 2011.

# Appendices

## A Proofs of the Theorems

### A.1 Proof of Theorem 1

*Proof.* The utility of a memory one player  $p$  against an opponent  $q$ ,  $u_q(p)$ , can be written as a ratio of two quadratic forms on  $R^4$ .

In Section 2, it was discussed that  $u_q(p)$  its the product of the steady states  $v$  and the PD payoffs,

$$u_q(p) = v \cdot (R, S, T, P).$$

More specifically,

$$u_q(p) = \frac{\begin{aligned} & p_1 p_2 (q_1 q_2 - 5q_1 q_4 - q_1 - q_2 q_3 + 5q_3 q_4 + q_3) + p_1 p_3 (-q_1 q_3 + q_2 q_3) + p_1 p_4 (5q_1 q_3 - 5q_3 q_4) + p_3 p_4 (-3q_2 q_3 + 3q_3 q_4) + \\ & p_2 p_3 (-q_1 q_2 + q_1 q_3 + 3q_2 q_4 + q_2 - 3q_3 q_4 - q_3) + p_2 p_4 (-5q_1 q_3 + 5q_1 q_4 + 3q_2 q_3 - 3q_2 q_4 + 2q_3 - 2q_4) + \\ & p_1 (-q_1 q_2 + 5q_1 q_4 + q_1) + p_2 (q_2 q_3 - q_2 - 5q_3 q_4 - q_3 + 5q_4 + 1) + p_3 (q_1 q_2 - q_2 q_3 - 3q_2 q_4 - q_2 + q_3) + \\ & p_4 (-5q_1 q_4 + 3q_2 q_4 + 5q_3 q_4 - 5q_3 + 2q_4) + q_2 - 5q_4 - 1 \end{aligned}}{\begin{aligned} & p_1 p_2 (q_1 q_2 - q_1 q_4 - q_1 - q_2 q_3 + q_3 q_4 + q_3) + p_1 p_3 (-q_1 q_3 + q_1 q_4 + q_2 q_3 - q_2 q_4) + p_1 p_4 (-q_1 q_2 + q_1 q_3 + q_1 + q_2 q_4 - q_3 q_4 - q_4) + \\ & p_2 p_3 (-q_1 q_2 + q_1 q_3 + q_2 q_4 + q_2 - q_3 q_4 - q_3) + p_2 p_4 (-q_1 q_3 + q_1 q_4 + q_2 q_3 - q_2 q_4) + p_3 p_4 (q_1 q_2 - q_1 q_4 - q_2 q_3 - q_2 + q_3 q_4 + q_4) + \\ & p_1 (-q_1 q_2 + q_1 q_4 + q_1) + p_2 (q_2 q_3 - q_2 - q_3 q_4 - q_3 + q_4 + 1) + p_3 (q_1 q_2 - q_2 q_3 - q_2 + q_3 - q_4) + p_4 (-q_1 q_4 + q_2 + q_3 q_4 - q_3 + q_4 - 1) + \\ & q_2 - q_4 - 1 \end{aligned}} \quad (23)$$

Let's consider the numerator of the  $u_q(p)$ . The cross product terms,  $p_i p_j$ , are given by

$$\begin{aligned} & p_1 p_2 (q_1 q_2 - 5q_1 q_4 - q_1 - q_2 q_3 + 5q_3 q_4 + q_3) + p_1 p_3 (-q_1 q_3 + q_2 q_3) + p_1 p_4 (5q_1 q_3 - 5q_3 q_4) + p_3 p_4 (-3q_2 q_3 + 3q_3 q_4) + \\ & p_2 p_3 (-q_1 q_2 + q_1 q_3 + 3q_2 q_4 + q_2 - 3q_3 q_4 - q_3) + p_2 p_4 (-5q_1 q_3 + 5q_1 q_4 + 3q_2 q_3 - 3q_2 q_4 + 2q_3 - 2q_4). \end{aligned}$$

The cross products' expression can be re written in a matrix format given by (24).

$$(p_1, p_2, p_3, p_4) \frac{1}{2} \begin{bmatrix} 0 & -(q_1 - q_3)(q_2 - 5q_4 - 1) & q_3(q_1 - q_2) & -5q_3(q_1 - q_4) \\ -(q_1 - q_3)(q_2 - 5q_4 - 1) & 0 & (q_2 - q_3)(q_1 - 3q_4 - 1) & (q_3 - q_4)(5q_1 - 3q_2 - 2) \\ q_3(q_1 - q_2) & (q_2 - q_3)(q_1 - 3q_4 - 1) & 0 & 3q_3(q_2 - q_4) \\ -5q_3(q_1 - q_4) & (q_3 - q_4)(5q_1 - 3q_2 - 2) & 3q_3(q_2 - q_4) & 0 \end{bmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} \quad (24)$$

Note that the coefficients are multiplied by  $\frac{1}{2}$  because they are added twice.

Similarly, the linear terms,

$$p_1(-q_1q_2 + 5q_1q_4 + q_1) + p_2(q_2q_3 - q_2 - 5q_3q_4 - q_3 + 5q_4 + 1) + p_3(q_1q_2 - q_2q_3 - 3q_2q_4 - q_2 + q_3) + p_4(-5q_1q_4 + 3q_2q_4 + 5q_3q_4 - 5q_3 + 2q_4).$$

can be written using a matrix format, (25).

$$(p_1, p_2, p_3, p_4) \begin{bmatrix} q_1(q_2 - 5q_4 - 1) \\ -(q_3 - 1)(q_2 - 5q_4 - 1) \\ -q_1q_2 + q_2q_3 + 3q_2q_4 + q_2 - q_3 \\ 5q_1q_4 - 3q_2q_4 - 5q_3q_4 + 5q_3 - 2q_4 \end{bmatrix} \quad (25)$$

Finally the constant term of the numerator, which is obtained by substituting  $p = (0, 0, 0, 0)$  is given by (26).

$$q_2 - 5q_4 - 1 \quad (26)$$

Equations (24), (25) and (26) are combined and the numerator of can be written as,

$$\frac{1}{2} p \begin{bmatrix} 0 & -(q_1 - q_3)(q_2 - 5q_4 - 1) & q_3(q_1 - q_2) & -5q_3(q_1 - q_4) \\ -(q_1 - q_3)(q_2 - 5q_4 - 1) & 0 & (q_2 - q_3)(q_1 - 3q_4 - 1) & (q_3 - q_4)(5q_1 - 3q_2 - 2) \\ q_3(q_1 - q_2) & (q_2 - q_3)(q_1 - 3q_4 - 1) & 0 & 3q_3(q_2 - q_4) \\ -5q_3(q_1 - q_4) & (q_3 - q_4)(5q_1 - 3q_2 - 2) & 3q_3(q_2 - q_4) & 0 \end{bmatrix} p^T +$$

$$\begin{bmatrix} 0 & -(q_1 - q_3)(q_2 - 5q_4 - 1) & q_3(q_1 - q_2) & -5q_3(q_1 - q_4) \\ -(q_1 - q_3)(q_2 - 5q_4 - 1) & 0 & (q_2 - q_3)(q_1 - 3q_4 - 1) & (q_3 - q_4)(5q_1 - 3q_2 - 2) \\ q_3(q_1 - q_2) & (q_2 - q_3)(q_1 - 3q_4 - 1) & 0 & 3q_3(q_2 - q_4) \\ -5q_3(q_1 - q_4) & (q_3 - q_4)(5q_1 - 3q_2 - 2) & 3q_3(q_2 - q_4) & 0 \end{bmatrix} p + q_2 - 5q_4 - 1$$

The same process is done for the denominator. □

## A.2 Proof of Theorem 3

*Proof.* Utility  $u_q(p)$  is non concave and neither are it's numerator or denominator.

A function  $f(x)$  is concave on an interval  $[a, b]$  if, for any two points  $x_1, x_2 \in [a, b]$  and any  $\lambda \in [0, 1]$ ,

$$f(\lambda x_1 + (1 - \lambda)x_2) \geq \lambda f(x_1) + (1 - \lambda)f(x_2). \quad (27)$$

Let  $f$  be  $u_{(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})}$ . For  $x_1 = (\frac{1}{4}, \frac{1}{2}, \frac{1}{5}, \frac{1}{2})$ ,  $x_2 = (\frac{8}{10}, \frac{1}{2}, \frac{9}{10}, \frac{7}{10})$  and  $\lambda = 0.1$ , direct substitution in (27) gives:

$$\begin{aligned} u_{(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})} \left( 0.1 \left( \frac{1}{4}, \frac{1}{2}, \frac{1}{5}, \frac{1}{2} \right) + 0.9 \left( \frac{8}{10}, \frac{1}{2}, \frac{9}{10}, \frac{7}{10} \right) \right) &\geq 0.1 \times u_{(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})} \left( \left( \frac{1}{4}, \frac{1}{2}, \frac{1}{5}, \frac{1}{2} \right) \right) + 0.9 \times u_{(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})} \left( \left( \frac{8}{10}, \frac{1}{2}, \frac{9}{10}, \frac{7}{10} \right) \right) \Rightarrow \\ 1.485 &\geq 0.1 \times 1.790 + 0.9 \times 1.457 \Rightarrow \\ 1.485 &\geq 1.490 \end{aligned}$$

which can not hold. Thus  $u_{(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})}$  is not concave. Because the concavity condition fails for at least one point of  $u_q(p)$ ,  $u_q(p)$  is not concave.

Utility  $u_q(p)$  is given by (6). As stated in [3] a quadratic form will be concave if and only if it's symmetric matrix is negative semi definite. A matrix  $A$  is semi-negative definite if:

$$|A|_i \leq 0 \text{ for } i \text{ is odd and } |A|_i \geq 0 \text{ for } i \text{ is even.} \quad (28)$$

For (6), neither  $\frac{1}{2}pQp^T + cp + a$  or  $\frac{1}{2}p\bar{Q}p^T + \bar{c}p + \bar{a}$  are concave because:

$$\begin{aligned} |Q|_2 &= -(q_1 - q_3)^2 (q_2 - 5q_4 - 1)^2 \text{ and} \\ |\bar{Q}|_2 &= -(q_1 - q_3)^2 (q_2 - q_4 - 1)^2 \end{aligned}$$

are negative. □