# When Reciprocity is not Enough: Explaining Anti-Social Outcomes via Intrinsic Personality Traits

PRELIMINARY VERSION

Discussion Paper

Michael Heinrich Baumann[*] and Michaela Baumann[†,‡]

June 24$^{th}$, 2025

*Abstract*—Behavior and outcomes that do not fit to classical game theory are often observed and, hence, reported. Common explanations for that are, e.g., repeated games, reciprocity, costly punishment, and pure altruism. Via the prisoner's dilemma and costly punishment, we show that those explanations, esp. reciprocity, are not always able to account for outcomes that involve anti-social punishment, i.e., to punish someone after a successful cooperation. We demonstrate, however, that intrinsic motivations, including both altruism and spiteful preferences, additional to material payoffs can explain those outcomes. To capture the agents' intrinsic motivation we introduce a uniform notion of altruism and sadism, the so-called SEA model. Further, we present a Python code to find so-called (pure-strategy) SEA Nash equilibria. Conclusively, we illustrate the SEA model via simple, well-known games.

[*]Department of Mathematics, University of Bayreuth, Germany, michael.baumann@uni-bayreuth.de

[†]NÜRNBERGER Versicherung, Nuremberg, Germany, michaela.baumann@nuernberger.de

[‡]Opinions expressed here are her own and not necessarily those of her employer.

1

# 1 Introduction and Literature Review

There are many situations in the real world as well as experiments that do not fit to predictions of classical game theory, for example, where the agents do not end up in a Nash equilibrium [11]. This phenomenon is well studied, e.g., on the basis of the prisoner's dilemma, see Table 1.

Table 1: Prisoner's Dilemma: material payoffs with values from [16]

| $u_1(\cdot)\|u_2(\cdot)$ | $a_2^{(1)}$ | $a_2^{(2)}$ |
|:---:|:---:|:---:|
| $a_1^{(1)}$ | 3\|3 | 0\|5 |
| $a_1^{(2)}$ | 5\|0 | 1\|1 |

The strategy (or action) $a_i^{(1)}$ refers to staying silent (or cooperation) and $a_i^{(2)}$ to confessing (or defecting). The outcome $(a_1^{(2)}, a_2^{(2)})$ is not only the only Nash equilibrium (cf. [11]), but also an equilibrium in strictly dominated strategies. As depicted in [3], in the references therein, prominently in [4, 5], agents in experiments and in the real world do (sometimes) cooperate, though.

A common explanation for cooperation in the prisoner's dilemma uses repeated games,[1] i.e. finitely or infinitely many repetitions of the very same game, which give the possibility to punish non-cooperation and reward co-operation—cf. the work on "tit for tat." The various other explanations that do also apply for one-shot games include pure altruism, a second round with another game that is designed especially for the possibility of punishment, and reciprocity. Concepts that use distributions/symmetry of payoffs among

---

[1]See, e.g., [16] Sections 3.4 and 3.5.

the agents do not directly apply in a meaningful (i.e., easy to interpret) way to the prisoner's dilemma since these are typically designed for $n$-player games (with large $n \in \mathbb{N}$).

In contrast to that, we introduce parameters for the agents reflecting their level of intrinsic altruism resp. intrinsic sadism without changing a game's objective (i.e. material) outcomes. Considering, however, the agents' subjective (i.e. psychological) outcomes influenced by these parameters, we are able to determine parameter combinations for which certain actions (or strategies) are Nash equilibria. For doing this, we firstly provide some short notes and references concerning reciprocity. After that, we introduce the concept of punishment, which brings us directly to the topic of so-called anti-social punishment, which is the main motivation for the work at hand. Thereupon, we are prepared to devise the so-called SEA model. We provide its definition, how it can be implemented, and give illustrative examples of its use.

## 1.1 Rabin's Fairness Equilibria

Rabin [13] uses the concept of reciprocity and beliefs, which he calls fairness, to explain outcomes that do not directly fit to Nash [11]. Applied on the prisoner's dilemma, this can be used to show that and why cooperation happens [1, 13]. This concept builds upon [8] and is based on three stylized facts, namely that agents agree to smaller payoff in order to be 1.) friendly to someone who is believed to be also friendly and 2.) unfriendly in order to hurt someone who is believed to be unfriendly, too. But, 3.) the larger those suffered fairness losses, the less do agents agree to those smaller payoffs.

In [1], the concept of Rabin fairness [13] is explained in great detail. There, a Python code is given, which is also used in the work at hand, to check whether outcomes are fair according to Rabin. For all details concerning Rabin fairness, please consult [1, 13].

To apply Rabin's fairness concept, we have to use a scaling factor $\chi > 0$, which accounts for the tradeoff between material payoff and fairness payoff, which is a payoff derived via so-called kindness functions. These kindness functions aim at measuring an agent's (un-)friendliness depending on the materials payoffs and agent's beliefs about the (un-)friendliness of the respective other.

The basic principle of fairness equilibria is the following: The (first-order believed) material payoff—which depends on the own action/strategy and the

3

first-order belief about what the opponent is going to do—is scaled by some $\chi > 0$ and there is a summand added, which depends on the (first-order believed) friendliness of the player towards the opponent—which depends also on the own action/strategy and the first-order belief about what the opponent is going to do—and on the (second-order believed) friendliness of the opponent towards the player—which depends on the first-order belief about what the opponent is going to do and on the second-order belief about what the opponent believes what the player is going to do. When these so-called kindness functions are either both positive or both negative, the fairness payoff is positive, when one is positive and one is negative, the fairness payoff is negative. In the equilibrium, the actions/strategies (and not the beliefs) have to maximize this so-called expected utility and actions and first- and second-order beliefs must match.

For the scaled prisoner's dilemma according to Rabin [13], see Table 2. One can show that (confess, confess) is for all scalings a fairness equilibrium while (stay silent, stay silent) is fair if and only if the suffered material losses for being friendly are not to high, in detail: for $\chi \leq 0.25$, see [1].

Table 2: Scaled Prisoner's Dilemma: material payoffs with values from [16] and scaled by $\chi > 0$ according to [13]

| $u_1(\cdot)\|u_2(\cdot)$ | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $3\chi\|3\chi$ | $0\|5\chi$ |
| $a_1^{(2)}$ | $5\chi\|0$ | $\chi\|\chi$ |

## 1.2 Costly Punishment

Another explanation for cooperation, which sounds quite meaningful in real world problems, is that there is a second round in an extended game called punishment. Often, this line of research does not directly use prisoner's dilemmas but so-called public goods games, which are quite similar to prisoner's dilemmas but allow for more agents. For example, in [12], a public goods game is utilized in which four agents can spend $a_i$ (money units) from zero up to 20 to the public and receive $u_i = (20 - a_i) + 0.4 \sum_j a_j$. Thus, if everyone gives 20, everyone gets 32. However, cheating (free riding) is strictly

dominant. The punishment round in [12] works as follows (if we understood it correctly): every agent was informed how much every other agent contributed and then could choose between zero and ten to punish one selected opponent with costs of factor one for the punisher and factor three for the punished.

Usually, stories to enhance plausibility of cooperation through costly punishment go as follows: Imagine a neighborhood where people should care about the common playground. Everyone profits when at least one cares about the playground, however, working there is hard and work is divided among volunteers, so, everyone wants to free ride and in the end no one works, which is not Pareto optimal. When there are barbecues in this neighborhood, but only volunteers are allowed to participate, there is some kind of costly punishment to free riders. Not being invited to the barbecue is very harmful for an individual. And also the volunteers profit if no one is excluded, since the fun is according to the-more-the-merrier. With this punishment (exclusion from the barbecue), everyone volunteers because he or she is afraid of not being invited. Note that the exclusion from the barbecue is something like an empty threat, since it is costly. Volunteering and punishing only the free riders is a Nash equilibrium (when formalized correctly). However, it is, due to the empty threat, not subgame perfect (see [16]). The punishment has to be costly since punishment-for-fun should be avoided (forcing egoistic agents not to punish, see Table 3).[2] That such an avoidance of punishment-for-fun does not work in all cases is exactly the topic of the main part of the work at hand.

There is a vast body of literature analyzing whether and to which extent costly punishment enhances cooperation (often, but not always, costly punishment enhances cooperation in the literature), see [14, 18, 19]. We explain the idea of costly punishment formally with the scaled prisoner's dilemma in Table 2 and the costly punishment of Table 3, again with a scaling factor $\varpi > 0$. Here, $a_i^{(3)}$ means no punishment of the other agent and $a_i^{(4)}$ means punishment of the other.

Clearly, if both agents play $a_i^{(2)}$ (confess) in the prisoner's dilemma and always no punishment in the second round $(a_i^{(3)}; a_i^{(3)}; a_i^{(3)}; a_i^{(3)})$, this is a Nash equilibrium, since $(a_1^{(2)}, a_2^{(2)})$ is an equilibrium in strictly dominated strategies in game 1 and $(a_1^{(3)}, a_2^{(3)})$ is an equilibrium in strictly dominated strategies in game 2. This equilibrium is subgame perfect. However, also if both agents

---

[2]Scaling is in the work at hand always done as in [13].

stay silent in the first game and punish the other one if and only if the other confessed in the first game, they play a Nash equilibrium (which is not subgame perfect since punishing is connected to a material loss of the punisher) if $\varpi$ (the punishment) is large enough compared to $\chi$ (the material payoff of game 1. If, e.g., $\chi = \varpi = 1$ no agent has an incentive to deviate: deviation in game 2 when the other cooperated would cause unnecessary costs, deviation in game 2 when the other confessed does not change the payoff since this is not played in the equilibrium by the other agent, and confessing would lead to a punishment that is higher than the payoff gained from cheating.

To make this two-round game easier to analyze, we could rewrite it into a one-shot game, cf. [10] Ch. 2.2 (and 2.6). This leads to a table where each agent got $2^5 = 32$ pure strategies, thus, in Table 4 we depicted only the structure. Here, we use the following notation (ordering):

- strategy in game 1;(

- strategy in game 2 if $(a_1^{(1)}, a_2^{(1)})$ was played in game 1;

- strategy in game 2 if $(a_1^{(2)}, a_2^{(1)})$ was played in game 1;

- strategy in game 2 if $(a_1^{(1)}, a_2^{(2)})$ was played in game 1;

- strategy in game 2 if $(a_1^{(2)}, a_2^{(2)})$ was played in game 1)

That is, we denote the actions for round 2 according to the outcomes of round 1 column-wise.

First, we calculate whether and when the strategy pair explained above, namely cooperation with the (empty) threat of punishing non-cooperative others (when used by both), is a fairness equilibrium or a Nash equilibrium (we already know that it is not a subgame perfect strategy). This is

Table 3: Scaled Costly Punishment: material payoffs scaled by $\varpi > 0$

| $u_1(\cdot)\|u_2(\cdot)$ | $a_2^{(3)}$ | $a_2^{(4)}$ |
|---|---|---|
| $a_1^{(3)}$ | $0\|0$ | $-10\varpi\|-\varpi$ |
| $a_1^{(4)}$ | $-\varpi\|-10\varpi$ | $-11\varpi\|-11\varpi$ |

6

Table 4: Scaled Prisoner's Dilemma with Scaled Costly Punishment: material payoffs with values from Tables 2 resp. [16] and 3 scaled by $\chi, \varpi > 0$ truncated

| $u_1(\cdot)|u_2(\cdot)$ | $a_2^{(1)}; (a_2^{(3)}; a_2^{(3)}; a_2^{(3)}; a_2^{(3)})$ | $\dots$ | $a_2^{(2)}; (a_2^{(4)}; a_2^{(4)}; a_2^{(4)}; a_2^{(4)})$ |
|---|---|---|---|
| $a_1^{(1)}; (a_1^{(3)}; a_1^{(3)}; a_1^{(3)}; a_1^{(3)})$ | $3\chi|3\chi$ | $\dots$ | $-10\varpi|5\chi-\varpi$ |
| $a_1^{(1)}; (a_1^{(3)}; a_1^{(3)}; a_1^{(3)}; a_1^{(4)})$ | $3\chi|3\chi$ | $\dots$ | $-10\varpi|5\chi-\varpi$ |
| $a_1^{(1)}; (a_1^{(3)}; a_1^{(3)}; a_1^{(4)}; a_1^{(3)})$ | $3\chi|3\chi$ | $\dots$ | $-11\varpi|5\chi-11\varpi$ |
| $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $a_1^{(2)}; (a_1^{(4)}; a_1^{(4)}; a_1^{(4)}; a_1^{(4)})$ | $5\chi-\varpi|-10\varpi$ | $\dots$ | $\chi-11\varpi|\chi-11\varpi$ |

the strategy pair which we call social punishment: $(a_1^{(1)}; (a_1^{(3)}; a_1^{(3)}; a_1^{(4)}; a_1^{(4)}),$ $a_2^{(1)}; (a_2^{(3)}; a_2^{(4)}; a_2^{(3)}; a_2^{(4)}))$

To enhance simplicity, we set $\varpi = \chi > 0$. Then we can use the Python [17] code (using SymPy [9]) from [1], which works for univariately scaled games only. It would be very interesting to analyze this two-rounds game with two (different) scaling variables (by calculations from hand or even via an adequate code), however, this lies beyond the scope of the work at hand and is postponed to future work. To analyze the prisoner's dilemma with punishment we insert in our code from [1] the new game.[3]

```
1  def prisoners_dilemma_punishment():
2      # prisoner's dilemma (Sieg)
3      v1 = [[3,0],[5,1]]
```

---

[3]Please consult Footnote 24 from [1], which we cite here for completeness: "For the inequality solver solve_poly_inequality see https://docs.sympy.org/latest/modules/solvers/inequalities.html (2024-03-21). For intervals, set operations, and oo ($\infty$), see https://docs.sympy.org/latest/modules/sets.html (2024-03-21). For the reduce function functools.reduce (fun,seq) see https://www.geeksforgeeks.org/reduce-in-python/ (2024-03-21). For the topics copy, deepcopy, and mutable objects, see https://stackoverflow.com/questions/8743072/when-adding-to-list-why-does-python-copy-values-instead-of-pointers (2024-03-25), https://stackoverflow.com/questions/19210971/python-prevent-copying-object-as-reference (2024-03-26), and https://docs.python.org/3/library/copy.html (2024-03-26). And, finally, for time, see https://www.python-lernen.de/python-modul-time.htm (in German; 2024-03-26)." See also https://www.geeksforgeeks.org/divmod-python-application/ (2025-05-27) for divmod.

```python
 4        v2 = copy.deepcopy(v1)
 5        # costly punishment (Baumann and Baumann)
 6        w1 = [[0,-10],[-1,-11]]
 7        w2 = copy.deepcopy(w1)
 8        u1 = []
 9        u2t = []
10
11        S1 = range(2**5)
12        S2 = range(2**5)
13        for i in S1:
14            u1.append([])
15            u2t.append([])
16            for j in S2:
17
18                a10,r = divmod(i,2**4)
19                a11,r = divmod(r,2**3)
20                a12,r = divmod(r,2**2)
21                a13,r = divmod(r,2)
22                a14 = r
23                a20,r = divmod(j,2**4)
24                a21,r = divmod(r,2**3)
25                a22,r = divmod(r,2**2)
26                a23,r = divmod(r,2)
27                a24 = r
28
29        # if ax0 is 0, agent x plays cooperate,
30        # if it's 1 he or she plays defect
31        # ax1 is the strategy in round 2 of agent x, when a10=0 and a20=0
32        # ax2 is the strategy in round 2 of agent x, when a10=1 and a20=0
33        # ax3 is the strategy in round 2 of agent x, when a10=0 and a20=1
34        # ax4 is the strategy in round 2 of agent x, when a10=1 and a20=1
35
36                u1[i].append(v1[a10][a20])
37                if a10 == 0 and a20 == 0: u1[i][j] = u1[i][j] + w1[a11][a21]
38                if a10 == 1 and a20 == 0: u1[i][j] = u1[i][j] + w1[a12][a22]
39                if a10 == 0 and a20 == 1: u1[i][j] = u1[i][j] + w1[a13][a23]
40                if a10 == 1 and a20 == 1: u1[i][j] = u1[i][j] + w1[a14][a24]
41
```

8

```
42          # be careful! u2 is transposed at the moment
43          u2t[i].append(v2[a20][a10])
44          if a10 == 0 and a20 == 0: u2t[i][j] = u2t[i][j] + w2[a21][a11]
45          if a10 == 1 and a20 == 0: u2t[i][j] = u2t[i][j] + w2[a22][a12]
46          if a10 == 0 and a20 == 1: u2t[i][j] = u2t[i][j] + w2[a23][a13]
47          if a10 == 1 and a20 == 1: u2t[i][j] = u2t[i][j] + w2[a24][a14]
48
49      # transposing u2
50      u2 = copy.deepcopy(u2t)
51      for i in S2:
52          for j in S1:
53              u2[i][j] = u2t[j][i]
54
55      return u1, u2
56      # end def prisoners_dilemma_punishment
```

It turns out that social punishment is a Nash equilibrium and that it is a fairness equilibrium for all $\chi > 0$. Playing Nash in all stages $(a_1^{(2)}; (a_1^{(3)}; a_1^{(3)}; a_1^{(3)}; a_1^{(3)}), a_2^{(2)}; (a_2^{(3)}; a_2^{(3)}; a_2^{(3)}; a_2^{(3)}))$ is a Nash equilibrium and a fairness equilibrium for all $\chi \geq \frac{5}{49}$. Thus, one could claim that social punishment is fairer than always playing Nash.

## 1.3   Anti-Social Punishment

Unintuitively, but nevertheless interestingly, when playing public goods games with costly punishment (which is nearly the same as the prisoner's dilemma with costly punishment) in experiments, so-called anti-social punishment occurs [12]. The authors of [12] analyze which factors correlate to the occurrence of anti-social punishment. That is, some agents are also punishing cooperative opponents.

Using our Python code [1] (with SymPy, cf. [9, 17]) we confirm that the corresponding strategies $(a_1^{(1)}; (a_1^{(4)}; a_1^{(4)}; a_1^{(4)}; a_1^{(4)}), a_2^{(1)}; (a_2^{(3)}; a_2^{(4)}; a_2^{(3)}; a_2^{(4)}))$ (if only one agent is anti-social punishing, i.e. punishing in all cases) and $(a_1^{(1)}; (a_1^{(4)}; a_1^{(4)}; a_1^{(4)}; a_1^{(4)}), a_2^{(1)}; (a_2^{(4)}; a_2^{(4)}; a_2^{(4)}; a_2^{(4)}))$ are neither Nash equilibria (which is clear) nor fairness equilibria (for all $\chi$), which is interesting since this shows that even a concept of reciprocity [13] resp. beliefs (that the other is not kind) cannot explain anti-social punishments.

9

## 1.4 Altruism

The last possibility to be discussed here for explaining non-Nash outcomes shall be altruism. That means, agents who have an intrinsic motivation to help others. For this concept please consult the review [7] and the references therein. This can be modeled in various ways, e.g., one could set

$$U_i(a_i, a_{-i}) := u_{-i}(a_{-i}, a_i),$$

which means that agents would care only about the other, or

$$U_i(a_i, a_{-i}) := u_i(a_i, a_{-i}) + u_{-i}(a_{-i}, a_i)$$

(or, equivalently multiplied by 0.5), which would mean that agents care equally for themselves and the other. We will use a generalization, which is inspired by [27],[4] but has the advantage that it is linear and, thus, easier to handle, which allows for continuously variable sizes of altruism, namely

$$U_i(\lambda_i, a_i, a_{-i}) := (1 - \lambda_i)u_i(a_i, a_{-i}) + \lambda_i u_{-i}(a_{-i}, a_i).$$

However, in the next section we are going to also incorporate the opposite of altruism, namely sadism (i.e. spiteful preferences). We do not analyze here whether pure altruism can explain anti-social punishment, since we will do a more general analysis in the next chapter.

## 2 Incorporating Intrinsic Altruism and Intrinsic Sadism

A possibility to explain one-sided anti-social punishment is by introducing intrinsic motivations other than material payoffs, namely intrinsic altruism and intrinsic sadism (i.e. personality traits). Altruism means that someone cares about the other, sadism means that someone wants to hurt the other.

---

[4]In [27], it is by means of the "battle of the sexes" (see also [13]) explained how non (standard) Nash outcomes can be explained, in order to adapt to real-world situations. For that, so-called Rawlsian functions are used: $U_i(\gamma_i, a_i, a_{-i}) := \min\{u_i(a_i, a_{-i}) \cdot \gamma_i^{-1}, u_{-i}(a_{-i}, a_i) \cdot (1 - \gamma_i)^{-1}\}$, cf. [15], which have the disadvantage to complicate calculations. This is why we propose linearly mixing in the altruism case.

## 2.1 The SEA Model

We construct a model that respecifies the outcomes according to the types of the agents, which are described with uniform parameters $\lambda_i \in [-1, 1]$ where $\lambda_i = -1$ is a pure sadist, $\lambda_i \in (-1, 0)$ is a sadist that also cares about him- or herself, $\lambda_i = 0$ is an egoist, $\lambda_i \in (0, 1)$ is an altruist that also cares about him- or herself, and $\lambda_i = 1$ is a pure altruist. These types are intrinsic in such a sense that agents are not kind to those who are kind to them, but they are always kind or not. The parameter $\lambda_i \in [-1, 1]$ allows for gradually shifting the type between the poles. We stick to common knowledge and to rationality. Also the types of the agents are common knowledge (which is in contrast to the experiments in [12], where the (types of the) opponents were not known, i.e., the agents had incomplete information). In the work at hand, we introduce this common knowledge model with complete information. In future works, this restriction may be loosened allowing for unknown agent types. Also, not least for the reason of avoiding circular reasoning, agents are one-step empathic, i.e., they care (positively or negatively) about the material payoff of the other, but not about his or her respecified payoff. Our respecified payoffs are continuous in the types and well-defined.

**Definition 1.** *We respecify the material payoffs $u_i(a_i, a_{-i})$ into psychological payoffs with $(\lambda_1, \lambda_2) \in [-1, 1]^2$ via*

$$U_i(\lambda_i, a_i, a_{-i}) = (1 - |\lambda_i|)u_i(a_i, a_{-i}) + \lambda_i u_{-i}(a_{-i}, a_i)$$

*for $i = 1, 2$. The game, when material payoffs are transformed via this formulae, is called the Sadism-Egoism-Altruism-model (SEA model).*

We highlight that we can gradually and continuously shift the types. Agents with $\lambda_i \in [-1, 0)$ can be called anti-social (although it is discussable whether a pure egoist is anti-social or not) and agents with $\lambda_i \in (0, 1]$ social. The larger $|\lambda_i|$ is, the more important is the material payoff of the opponent. The larger $\lambda_i$ is, the more altruistic is agent $i$, the smaller $\lambda_i$ gets, the more he or she is sadisdic, and if it is zero, he or she is an egoist. If both $\lambda_1$ and $\lambda_2$ are zero, the game is the original game.

We define an outcome (i.e. a pair of strategies) as

- SEA Nash equilibrium for $\mathbf{E} \subset [-1, 1]^2$ if the outcome is a Nash equilibrium [11] in the SEA model for all $(\lambda_1, \lambda_2) \in \mathbf{E}$ and not a Nash equilibrium in the SEA model for all $(\lambda_1, \lambda_2) \in [-1, 1]^2 \setminus \mathbf{E}$,

11

- SEA plausible, if there exists a $(\lambda_1, \lambda_2) \in [-1, 1]^2$ such that the outcome is a Nash equilibrium in the SEA model with $\lambda_1, \lambda_2$,

- SEA plausible under social types, if there exists a $(\lambda_1, \lambda_2) \in ([0, 1] \times (0, 1]) \cup ((0, 1] \times [0, 1]) = [0, 1]^2 \setminus \{(0, 0)\}$ such that the outcome is a Nash equilibrium in the SEA model with $\lambda_1, \lambda_2$,

- SEA plausible under anti-social types, if there exists a $(\lambda_1, \lambda_2) \in ([-1, 0] \times [-1, 0)) \cup ([-1, 0) \times [-1, 0]) = [-1, 0]^2 \setminus \{(0, 0)\}$ such that the outcome is a Nash equilibrium in the SEA model with $\lambda_1, \lambda_2$,

- SEA plausible under mixed types, if there exists a $(\lambda_1, \lambda_2) \in ([-1, 0) \times (0, 1]) \cup ((0, 1] \times [-1, 0))$ such that the outcome is a Nash equilibrium in the SEA model with $\lambda_1, \lambda_2$.

We note that this model generalizes in some sense the idea of reciprocity: Being kind to those who are (supposed to be) kind to you corresponds to the concept of SEA plausibility under social types. Wanting someone hurt who wants you hurt (or is supposed to do so) corresponds to SEA plausibility under anti-social types. Further, the SEA model allows for the asymmetric cases, however, according to the signs of the $\lambda_i$, we can distinguish these cases (also the two different reciprocal cases). We note that at first, we stick to pure strategies and outcomes only, however, the generalization to mixed strategies is straightforward. That every Nash equilibrium is SEA plausible follows by setting $\lambda_1 = \lambda_2 = 0$:

**Proposition 1.** *Every Nash equilibrium (of the material game) is SEA plausible.*

For interpretations it is important to remark that agents know whether and to what degree the opponent is sadistic or altruistic and they behave nevertheless according to their type. Some help their enemies because they want to or it is their nature to do so and some bite the hand who feeds them because it is their nature resp. because they want to do so. Thus, results do not fit to the experiment setting of [12], however, they may explain why and when anti-social punishment happens.

We emphasize that in general not every outcome is SEA plausible and if so, this does not have to hold for all parameter combinations.

## 2.2 The Shape of the SEA Nash Sets

Since we aim at providing a Python code to automatically compute the SEA Nash equilibria, we first provide a proposition on the shape of the sets $\mathbf{E}$. This information is very helpful when implementing the code.

**Proposition 2.** *The sets $\mathbf{E} \subset [-1,1]^2$ in the definition of SEA Nash equilibria are always rectangles with sides parallel to the axes.*

*Proof.* Let the outcome $(a_1^\star)$ be SEA Nash for $\mathbf{E} \ni (\lambda_1, \lambda_2), (\lambda_1', \lambda_2')$. We will show firstly that $(\lambda_1, \lambda_2') \in \mathbf{E}$.

For that, $(\lambda_1, \lambda_2) \in \mathbf{E}$ means that

$$U_1(\lambda_1, a_1^\star, a_2^\star) \geq U_1(\lambda_1, a_1, a_2^\star) \ \forall a_1 \in A_1$$

and $U_2(\lambda_2, \ldots,$ cf. [11]. Further, $(\lambda_1', \lambda_2') \in \mathbf{E}$ means that

$$U_2(\lambda_2', a_2^\star, a_1^\star) \geq U_2(\lambda_2', a_2, a_1^\star) \ \forall a_2 \in A_2$$

and $U_1(\lambda_1', \ldots.$

These two inequalities show that $(a_1^\star, a_2^\star)$ is SEA Nash for $(\lambda_1, \lambda_2')$. Analogously holds $(\lambda_1', \lambda_2) \in \mathbf{E}$, which completes the proof.

$\square$

## 2.3 Python/SymPy Code

In this section we provide some Python code [9, 17], which can be inserted in the code provided in [1]. Additionally to the Rabin fairness equilibria, Nash equilibria, and Mutual Min resp. Mutual Max outcomes we can then also calculate the SEA Nash equilibria.

```
57  def sea_nash_equilibria(game):
58
59      ### The SEA-Model (Sadism-Egoism-Altruism)
60
61      # Intrinsic Altruism and Intrinsic Sadism
62
63      # Ordering of list elements according to lambda_1=y, lambda_2=z:
64      # [y,z >= 0; y >= 0 and z < 0; y < 0 and z >= 0; y,z < 0]
65
```

13

```python
66        t = time.time()
67
68        # material payoffs
69        u1, u2 = set_game(game)
70
71        # Altruism and Sadism Parameters for agent 1 (y) and agent 2 (z)
72        y = sympy.symbols('y')
73        z = sympy.symbols('z')
74
75        # actions
76        S1 = range(len(u1))
77        S2 = range(len(u1[0]))
78
79        # psychological payoffs according to Definition 1
80        e1 = []
81        e2 = []
82        for l in range(4):
83            e1.append(copy.deepcopy(u1))
84            e2.append(copy.deepcopy(u2))
85
86        for i in S1:
87            for j in S2:
88                e1[0][i][j] = (1-y)*u1[i][j]+y*u2[j][i]
89                e1[1][i][j] = (1-y)*u1[i][j]+y*u2[j][i]
90                e1[2][i][j] = (1+y)*u1[i][j]+y*u2[j][i]
91                e1[3][i][j] = (1+y)*u1[i][j]+y*u2[j][i]
92                e2[0][j][i] = (1-z)*u2[j][i]+z*u1[i][j]
93                e2[1][j][i] = (1+z)*u2[j][i]+z*u1[i][j]
94                e2[2][j][i] = (1-z)*u2[j][i]+z*u1[i][j]
95                e2[3][j][i] = (1+z)*u2[j][i]+z*u1[i][j]
96
97
98        # Nash equilibria
99        # SEA best responses
100        SeaBR1 = [[],[],[],[]]
101        SeaBR2 = [[],[],[],[]]
102        SeaNash = [[],[],[],[]]
103        for l in range(4):
```

14

```
104            # best response functions
105
106            # initial values
107
108        for i in S1:
109            SeaBR1[l].append([])
110            SeaBR2[l].append([])
111            SeaNash[l].append([])
112            for j in S2:
113                if l == 0 or l == 1:
114                    SeaBR1[l][i].append(Interval(0,1))
115                else:
116                    SeaBR1[l][i].append(Interval.Ropen(-1,0))
117                if l == 0 or l == 2:
118                    SeaBR2[l][i].append(Interval(0,1))
119                else:
120                    SeaBR2[l][i].append(Interval.Ropen(-1,0))
121                SeaNash[l][i].append(sympy.EmptySet)
122
123                # is i in S1 a best response if agent 2 plays j in S2?
124
125                for k in S1:
126
127                    SeaBR1[l][i][j] = functools.reduce(
128                        lambda a, b: a.union(b), (
129                        sympy.solve_poly_inequality(
130                        Poly(e1[l][k][j]-e1[l][i][j],y,domain='RR'), ">")
131                        )
132                        ).complement(Interval(-1,1)
133                        ).intersection(SeaBR1[l][i][j])
134
135                for k in S2:
136
137                    SeaBR2[l][i][j] = functools.reduce(
138                        lambda a, b: a.union(b), (
139                        sympy.solve_poly_inequality(
140                        Poly(e2[l][k][i]-e2[l][j][i],z,domain='RR'), ">")
141                        )
```

```
142                    ).complement(Interval(-1,1)
143                    ).intersection(SeaBR2[l][i][j])
144
145              SeaNash[l][i][j]=[SeaBR1[l][i][j],SeaBR2[l][i][j]]
146
147        # for interpretation, consult what is written in Baumann and Baumann
148        # 2025 Some Thoughts on Rabin Fairness:
149        # "BR1 = [[0, 0, 1], [1, 0, 0], [0, 1, 0]] means that i's first
150        # strategy is the best response to -i's third one, i's second one is
151        # the best response to -i's first one, and finally i's third one is
152        # the best answer to -i's second strategy
153        #
154        # only pure and no mixed strategies and Nash equilibria are
155        # considered
156        #
157        # BR2 = [[0, 1, 0], [0, 0, 1], [1, 0, 0]] means that the best -i
158        # can do if i does its 1st, is its 2nd, the best -i can do if i
159        # plays its 2nd, is its 3rd, ..."
160
161        # union over all four cases
162
163    SeaNashUnion = copy.deepcopy(u1)
164    for i in S1:
165        for j in S2:
166            SeaNashUnion[i][j] = [SeaNash[0][i][j][0].union(SeaNash[1][i][j][0]
167            ).union(SeaNash[2][i][j][0]
168            ).union(SeaNash[3][i][j][0]),
169            SeaNash[0][i][j][1].union(SeaNash[1][i][j][1]
170            ).union(SeaNash[2][i][j][1]
171            ).union(SeaNash[3][i][j][1])
172            ]
173
174    runtime = time.time()-t
175
176    return SeaNashUnion, runtime
177    # end def sea_nash_equilibria
```

This code can be applied to finite two-agent games in normal form to

calculate pure-strategy SEA Nash equilibria. However, please note that we can neither prove the correctness of the code nor guarantee for it. First, we will apply it to the prisoner's dilemma with costly punishment.[5]

# 3 SEA Nash Equilibria for the Prisoner's Dilemma with Costly Punishment

For our prisoner's dilemma with costly punishment (unscaled, i.e., $\chi = \varpi = 1$) it turns out that playing Nash in all stages is SEA Nash for $\left[-\frac{1}{11}, \frac{1}{5}\right] \times \left[-\frac{1}{11}, \frac{1}{5}\right]$. Social punishment is SEA Nash for $\left[-\frac{1}{11}, 1\right] \times \left[-\frac{1}{11}, 1\right]$. Anti-social punishment (when agent 1 is punishing all the time and agent two only in the defect cases) is SEA Nash for $\left[-\frac{2}{3}, -\frac{1}{11}\right] \times \left[\frac{2}{5}, 1\right]$ and that both agents are punishing all the time (but cooperate in round 1) is not SEA plausible. These computations were done using the Python code from [1] and its extension depicted in Section 2.3 (cf. [9, 17]). The most important results are formulated in the next proposition again.

**Proposition 3.** *In the prisoner's dilemma with punishment as specified in Tables 2 and 3 with $\chi = \varpi = 1$ a one-sided anti-social punishment as defined above is neither Nash nor fair (in the sense of Rabin; for no $\chi > 0$ it is fair). However, it is SEA plausible, specifically it is SEA Nash for $\left[-\frac{2}{3}, -\frac{1}{11}\right] \times \left[\frac{2}{5}, 1\right]$.*

One-sided anti-social punishment is SEA plausible under mixed types, esp. the anti-social punisher has to be sadistic (but not too strong—in Round 1 we need cooperation) and the social punisher has to be altruistic (so he or she is not cheating in Round 1). We highlight that this one-sided anti-social punishment is SEA plausible only under mixed types. That means, one agent has to be sadistic and one has to be altruistic (each to some degree). This is interesting because altruism in games is widely discussed and analyzed in the literature [7] both as reciprocal altruism and as intrinsic altruism (in the literature also called pure altruism) while the body of literature dealing with sadism or spiteful preferences in games (resp. in game theory) is relatively

---

[5]Please mind the distinction between hurting and punishing: Hurting means to be unkind/unfriendly/mean to someone who is (believed to be) unkind/unfriendly/mean to you in this game; punishing means to act in a second game or future round of the same game in such a way that the opponent has some losses if he or she does not behave kind now. Punishment is usually used as an (possibly empty) threat.

small and often empirical [12]. Next, the code for displaying the discussed outcomes is provided.

```python
def print_specific_results(nash, fair, sea_nash):
    print("Playing Nash in all stages")
    i = 1*16+0*8+0*4+0*2+0
    j = 1*16+0*8+0*4+0*2+0
    print(nash[i][j])
    print(fair[i][j])
    print(sea_nash[i][j])

    print("Social punishment")
    i = 0*16+0*8+0*4+1*2+1
    j = 0*16+0*8+1*4+0*2+1
    print(nash[i][j])
    print(fair[i][j])
    print(sea_nash[i][j])

    print("One Social and One Anti-Social Punishment")
    i = 0*16+1*8+1*4+1*2+1
    j = 0*16+0*8+1*4+0*2+1
    print(nash[i][j])
    print(fair[i][j])
    print(sea_nash[i][j])

    print("One Social and One Anti-Social Punishment (check)")
    i = 0*16+0*8+0*4+1*2+1
    j = 0*16+1*8+1*4+1*2+1
    print(nash[i][j])
    print(fair[i][j])
    print(sea_nash[i][j])

    print("Two Anti-Social Punishments")
    i = 0*16+1*8+1*4+1*2+1
    j = 0*16+1*8+1*4+1*2+1
    print(nash[i][j])
    print(fair[i][j])
    print(sea_nash[i][j])
```

18

```
214       # end def print_specific_results
```

We highlight that the outcomes that two sadists cooperate in the prisoners dilemma, i.e. cooperation in Round 1 and always punish in Round 2, is not SEA plausible. Likely, this is due to the fact that our game is modeled with complete information, while the experiment in [12] is conducted with incomplete information. In an experiment it cannot be excluded that two sadists play with each other not knowing that the respective opponent is sadistic, too.

# 4   Simulation-based Comparison of Nash, Rabin, and the SEA Model in Other Games

So far, we have seen that the SEA model is powerful enough to explain (one-sided) anti-social punishment. In order to illustrate the SEA model further, we provide a couple of more or less classical and simple games hereafter and calculate (using our code) the Nash equilibria [11] in pure strategies, the fairness equilibria [13], and the SEA Nash equilibria under pure strategies. We stick to the standard form of games: higher values are favorable for the agents, there are two agents and each one has a finite set of actions leading to bounded outcomes. For the analyzed games we provide the payoff bi-matrix, a matrix with zeros and ones indicating which outcomes are Nash equilibria, a matrix with components that are subsets of $(0, \infty)$, indicating for which $\chi > 0$ the respective outcome is a fairness equilibrium [1, 13], and a matrix with components that are subsets of $[-1, 1]^2$, indicating for which $(\lambda_1, \lambda_2)$ the outcome is SEA Nash.

## 4.1   Prisoner's Dilemma (Sieg)

We start with our prisoner's dilemma from Table 1 with values from [16]. In Tables 5, 6, and 7 the results are depicted. Only (defect,defect) is Nash, which is also fair for all scaling factors. For small enough scaling factors, (cooperation,cooperation) is also fair. One-sided cooperation is never fair. If both agents are altruistic enough $(\lambda_1, \lambda_2 \geq 0.4)$ (cooperation,cooperation) is SEA Nash, if both are sadistic or at least not too altruistic $(\lambda_1, \lambda_2 \leq 0.2)$ (defect,defect) is SEA Nash. Interestingly, also one-sided cooperation is SEA

19

Nash—and this is not only the case when the cooperator is altruistic and the defector is sadistic, but also when the cooperator is less altruistic than the defector, e.g., $(\lambda_1, \lambda_2) = (0.2, 0.4)$.

Table 5: Prisoner's Dilemma with values from [16] or [13]: Nash equilibrium

| Nash | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | 0 | 0 |
| $a_1^{(2)}$ | 0 | 1 |

Table 6: Prisoner's Dilemma with values from [16] or [13]: Fairness equilibria

| Rabin | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $(0, 0.25]$ | $\emptyset$ |
| $a_1^{(2)}$ | $\emptyset$ | $(0, \infty)$ |

Table 7: Prisoner's Dilemma with values from [16]: SEA Nash equilibria

| SEA | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $[0.4, 1]^2$ | $[0.2, 1] \times [-1, 0.4]$ |
| $a_1^{(2)}$ | $[-1, 0.4] \times [0.2, 1]$ | $[-1, 0.2]^2$ |

## 4.2 Prisoner's Dilemma (Rabin)

As mentioned in [1], the exact values in the prisoner's dilemmas of Sieg [16] and Rabin [13] are slightly different while the structure—of course—is the same. Rabin uses a scaled version of Table 8. We observe that both the Nash equilibria and the fairness equilibria (Rabin) do not change. In Table 9 we see that the exact boundaries of the intervals when using the SEA model change, however, the structure does not.

Table 8: Prisoner's Dilemma: material payoffs with values from [13] with scaling parameter equal to one)

| $u_1(\cdot); u_2(\cdot)$ | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $4; 4$ | $0; 6$ |
| $a_1^{(2)}$ | $6; 0$ | $1; 1$ |

Table 9: Prisoner's Dilemma(with values from [13]): SEA Nash equilibria

| SEA | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $[0.\bar{3}, 1]^2$ | $[0.1\bar{6}, 1] \times [-1, 0.\bar{3}]$ |
| $a_1^{(2)}$ | $[-1, 0.\bar{3}] \times [0.1\bar{6}, 1]$ | $[-1, 0.1\bar{6}]^2$ |

## 4.3 Rock-Scissors-Paper (Sieg)

An important classical game in game theory is rock-scissors-paper, which is formalized by Sieg [16] as in Table 10. There is neither a Nash nor a fairness equilibrium in pure strategies, cf. [1], see Tables 11 and 12.

Table 10: Rock-Scissors-Paper with values from [16]: material payoffs

| $u_1(\cdot); u_2(\cdot)$ | $a_2^{(1)}$ | $a_2^{(2)}$ | $a_2^{(3)}$ |
|---|---|---|---|
| $a_1^{(1)}$ | $0; 0$ | $1; -1$ | $-1; 1$ |
| $a_1^{(2)}$ | $-1; 1$ | $0; 0$ | $1; -1$ |
| $a_1^{(1)}$ | $1; -1$ | $-1; 1$ | $0; 0$ |

Because all outcomes are SEA plausible, the insights are not deep, see Table 13. Despite the fact that rock-scissors-paper indicated that $\lambda_i = 0.5$ is an important threshold in the SEA model, we observe that agents agree in letting one win and one loose if one agent is at least one half altruistic and the other one is not more than one half altruistic or even (to any degree) sadistic.

Table 11: Rock-Scissors-Paper with values from [16]: Nash equilibria

| Nash | $a_2^{(1)}$ | $a_2^{(2)}$ | $a_2^{(3)}$ |
|------|------|------|------|
| $a_1^{(1)}$ | 0 | 0 | 0 |
| $a_1^{(2)}$ | 0 | 0 | 0 |
| $a_1^{(3)}$ | 0 | 0 | 0 |

Table 12: Rock-Scissors-Paper with values from [16]: fairness equilibria

| Rabin | $a_2^{(1)}$ | $a_2^{(2)}$ | $a_2^{(3)}$ |
|-------|------|------|------|
| $a_1^{(1)}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $a_1^{(2)}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |
| $a_1^{(3)}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |

## 4.4   Chicken (Rabin)

Last, we are going to analyze the chicken game, please consult [1, 13]. It can be formalized as done in [13], see Table 14.

The chicken game is particularly interesting for several reasons. 1.) Such as in [1, 13], we observe that there is no pure-strategy fairness equilibrium that is one for all $\chi > 0$, see Table 16. 2.) In [13], there is a distinction between strictly positive fairness equilibria (where both agents behave kindly) and weakly negative ones (where both agents do not behave kindly).[6] By, e.g., using the code from [1] it turns out that for $(a_1^{(1)}, a_2^{(1)})$, i.e. $(dare, dare)$ both kindness values are $-1$, which means that both agents are as mean as possible. For $(a_1^{(2)}, a_2^{(2)})$, i.e. $(chickenout, chickenout)$ both kindness values are 0.5, i.e. that both agents are as kind as possible. And for the outcomes where one agent "chickens out" and one "dares," that one who dares is "rationally mean" (i.e. 0.5; cf. [1]) and the "chicken" agent is neither kind nor mean (0). Those lastly mentioned two outcomes are Nash, see Table 15, the other two are not. 3.) When having a look at Table 17 we can gain much more insight. Really, $(dare, dare)$ is only SEA plausible when both are unkind

---

[6]See [13], esp. Definitions 1, 2, and 6 as well as Proposition 2.

Table 13: Rock-Scissors-Paper with values from [16]: SEA Nash equilibria

| SEA | $a_2^{(1)}$ | $a_2^{(2)}$ | $a_2^{(3)}$ |
|---|---|---|---|
| $a_1^{(1)}$ | $\{0.5\}^2$ | $[-1, 0.5] \times [0.5, 1]$ | $[0.5, 1] \times [-1, 0.5]$ |
| $a_1^{(2)}$ | $[0.5, 1] \times [-1, 0.5]$ | $\{0.5\}^2$ | $[-1, 0.5] \times [0.5, 1]$ |
| $a_1^{(3)}$ | $[-1, 0.5] \times [0.5, 1]$ | $[0.5, 1] \times [-1, 0.5]$ | $\{0.5\}^2$ |

Table 14: Chicken with values from [13]: material payoffs

| $u_1(\cdot); u_2(\cdot)$ | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $-2; -2$ | $2; 0$ |
| $a_1^{(2)}$ | $0; 2$ | $1; 1$ |

(sadisdic), and (*chickenout*, *chickenout*) only when both are kind (altruistic). However, the (*dare*, *chickenout*) and (*chickenout*, *dare*) outcomes are SEA plausible for all combinations of (not too) sadistic, egoistic, and (not too) altruistic agents.

# 5 Conclusion and Future Work

We showed that anti-social punishment can neither be explained by the concept of Nash [11] nor by that of Rabin (fairness [13]). However, incorporating intrinsic sadism and intrinsic altruism into the game can explain (one-sided) anti-social punishment. For this we introduced the SEA (Sadism-Egoism-Altruism) model, which allows to alter agent types gradually and continuously from pure sadists via egoists to pure altruists in a one-step empathic sense. This concept comes without beliefs. We give a Python code to compute those SEA Nash equilibria. Via simple and classical games, the SEA model is illustrated.

We emphasize that the SEA model generalizes in some sense the concept of reciprocity, where mutually "being kind" mirrors SEA plausibility under social types and mutually "being mean" mirrors SEA plausibility under anti-

Table 15: Chicken with values from [13]: Nash equilibria

| Nash | $a_2^{(1)}$ | $a_2^{(2)}$ |
|------|------|------|
| $a_1^{(1)}$ | 0 | 1 |
| $a_1^{(2)}$ | 1 | 0 |

Table 16: Chickenwith values from [13]: fairness equilibria

| Rabin | $a_2^{(1)}$ | $a_2^{(2)}$ |
|------|------|------|
| $a_1^{(1)}$ | $(0, 0.5]$ | $[0.25, \infty)$ |
| $a_1^{(2)}$ | $[0.25, \infty)$ | $(0, 0.5]$ |

social types. By means of the chicken game this is illustrated and it is shown that the SEA types can give more insights than the kindness values of [13]. Further, in the SEA model, plausibility under mixed types is possible.

Additional to the fairness concept of Rabin [13], there are also common and popular fairness concepts from Bolton and Ockenfels [2] and Fehr and Schmidt [6]. It is important for future work to check whether the strategies of interest (Nash in all stages, social punishment, one-sided anti-social punishment, mutual anti-social punishment) are fair in the senses of Bolton and Ockenfels as well as Fehr and Schmidt. Connections to the Minmax-theory of von Neumann (and Morgenstern) [20, 22, 21] and the concepts used by Wald [23, 24, 25, 26] and the mutual-min resp. mutual-max definitions of Rabin [13] are interesting, too.

The SEA model should be investigated thoroughly in theory and application. Impacts on economic questions and policy recommendations have to be analyzed as well. Maybe, the SEA model can be merged with Rabin fairness to explain social behavior even better. Additionally, the possibility of "agents doing failures" and the connected risk may be incorporated. And evolutionary aspects concerning sadism and altruism can be analyzed for or by means of the SEA model.

Compared to Rabin Fairness [13], the SEA model got the advantage that it is much more easy to compute, both by hand and also via our codes (see

24

Table 17: Chicken with values from [13]: SEA Nash equilibria

| SEA | $a_2^{(1)}$ | $a_2^{(2)}$ |
|---|---|---|
| $a_1^{(1)}$ | $[-1, -0.\bar{3}]^2$ | $[-1, 0.5] \times [-0.\bar{3}, 1]$ |
| $a_1^{(2)}$ | $[-0.\bar{3}, 1] \times [-1, 0.5]$ | $[0.5, 1]^2$ |

above and [1]). Rabin's concept has the advantage that the thinking in beliefs might better describe how real humans come to their decision, see [1]. However, opposite to the last topic, for analyzing behavior the SEA model can be favorable for scientists resp. researches.

# Acknowledgment

# References

[1] Baumann, Michael Heinrich, Michaela Baumann: Some Thoughts on Rabin Fairness. *University of Bayreuth,* Discussion Paper, 2025. `https://doi.org/10.15495/EPub_UBT_00008439`

[2] Bolton, Gary E., Axel Ockenfels: ERC: A Theory of Equity, Reciprocity, and Competition. *American Economic Review,* **91**(1):166-193, 2000.

[3] Cooper, Russell, Douglas V. DeJong, Robert Forsythe, Thomas W. Ross: Cooperation Without Reputation: Experimental Evidence From Prisoner's Dilemma Games. *Games and Economic Behavior,* **12**(2): 187-218, 1996.

[4] Dawes, Robyn Mason: Social Dilemmas. *Annual Review of Psychology,* **31**(1): 169-193, 1980.

[5] Dawes, Robyn Mason, Richard H. Thaler: Anomalies: Cooperation. *Journal of Economic Perspectives,* **2**(3): 187-197, 1988.

[6] Fehr, Ernst, Klaus M. Schmidt: A Theory of Fairness, Competition, and Cooperation. *The Quarterly Journal of Economics,* **144**(3):817-868, 1999.

[7] Fehr, Ernst, Klaus M. Schmidt: The Economics of Fairness, Reciprocity and Altruism – Experimental Evidence and New Theories. *In: Serge-Christophe Kolm, Jean Mercier Ythier (Eds.), Handbook of the Economics of Giving, Altruism and Reciprocity,* **1**:615-691, 2006.

[8] Geanakoplos, John, David Pearce, Ennio Stacchetti: Psychological Games and Sequential Rationality. *Games and Economic Behavior,* **1**: 60-79, March 1989.

[9] Meurer, Aaron, Christopher P. Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B. Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, Anthony Scopatz: SymPy: Symbolic Computing in Python. *Python, Computer algebra system, Symbolics,* **3**:e103, January 2017.

[10] Myerson, Roger B.: *Game Theory : Analysis of Conflict,* Harvard University Press, Cambridge, Massachusetts/London, England, 1991.

[11] Nash, John Forbes Jr.: Non-Cooperative Games. *Dissertation,* Princeton University, 1950.

[12] Pfattheicher, Stefan, Johannes Keller, Goran Knezevic: Sadism, the Intuitive System, and Antisocial Punishment in the Public Goods Game. *Personality and Social Psychology Bulletin,* **43**(3):337-346, 2017.

[13] Rabin, Matthew: Incorporating Fairness into Game Theory and Economics. *The American Economic Review,* **83**(5): 1281-1302, December 1993.

[14] Roberts, Gilbert: When Punishment Pays. *PLOS ONE,* **8**(3)e57378: 1-8, March 2013.

[15] Rawls, John: *A Theory of Justice,* Harvard University Press, Cambridge, Massachusetts, 1971.

[16] Sieg, Gernot: *Spieltheorie,* 2nd Edition, R. Oldenburg Verlag, München/Wien, 2005. *(in German)*

[17] Van Rossum, Guido, Fred L. Drake Jr.: *Python Tutorial.* Centrum voor Wiskunde en Informatica Amsterdam, The Netherlands, 1995

[18] Wu, Jia-Jia, Bo-Yu Zhang, Zhen-Xing Zhou, Qiao-Qiao He, Xiu-Deng Zheng, Ross Cressman, and Yi Tao: Costly Punishment does not Always Increase Cooperation. *Proceedings of the National Academy of Sciences,* **106**(41):17448-17451, October 2009.

[19] Ye, Hang, Fei Tan, Mei Ding, Yongmin Jia, Yefeng Chen: Sympathy and Punishment: Evolution of Cooperation in Public Goods Game. *Journal of Artificial Societies and Social Simulation,* **14**(4)20:1-14, 2011.

[20] von Neumann, John: Zur Theorie der Gesellschaftsspiele, *Mathematische Annalen,* 100:295-320, December 1928. *(in German)*

[21] von Neumann, John, Oskar Morgenstern: *Theory of Games and Economic Behaviour,* Princeton University Press, 1944.

[22] von Neumann, John: Über ein ökonomisches Gleichungssystem und eine Verallgemeinerung des Brouwerschen Fixpunktsatzes. *Ergebnisse eines Mathematischen Kolloquiums,* **8**(1935/36), Leipzig, 1937. *(in German)*

[23] Wald, Abraham: Contributions to the Theory of Statistical Estimation and Testing Hypotheses. *The Annals of Mathematics,* **10**(4):299-326, 1939.

[24] Wald, Abraham: Statistical Decision Functions which Minimize the Maximum Risk. *The Annals of Mathematics,* **46**(2):265-280, 1945.

[25] Wald, Abraham: *Statistical Decision Functions,* John Wiley, New York, 1950.

[26] Wald, Abraham: Generalization of a Theorem of von Neumann Concerning Zero Sum Two Person Games. *The Annals of Mathematics,* **46**(2):281-286, 1945.

[27] Zapata, Asunción , Amparo M. Mármol, Luisa Monroy, M. Ángeles Caraballo: When the Other Matters. The Battle of the Sexes Revisited. In: Patrizia Daniele, Laura Scrimali (Eds.) *New Trends in Emerging Complex Real Life Problems,* 501-509, AIRO Springer Series 1, ODS, Taormina, Italy, September 2018.