

# Social Influence on Third-Party Punishment: an Experiment<sup>☆</sup>

Emanuela Carbonara<sup>a</sup>, Marco Fabbri<sup>b</sup>

<sup>a</sup>*School of Economics, University of Bologna*

<sup>b</sup>*Institute of Private Law, Erasmus University Rotterdam*

---

## Abstract

In this paper we study the effects of social influence on third-parties' decision to engage in decentralized costly punishment. We elicit punishment decisions both in isolation and providing information about actual peers' punishment. We find evidence that that the amount of punishment chosen by third-parties is influenced by beliefs about the amount of peers' punishment. Moreover, the larger the difference between third-parties beliefs about the level of peers' punishment and actual peers' punishment, the more likely it is that third-parties modify their initial punishment decision. We also find that more self-regarding third-parties are less affected by social influence. Finally, we disentangle the effect of *Normative* social influence from that of *Informational* social influence and we show that the former type of social influence may be effective on subjects tending to disregard the latter.

*Keywords:* Third-party Punishment, Social Influence, Peer Pressure, Experiment, Dictator Game

---

<sup>☆</sup>Contact author: Marco Fabbri, [fabbri@law.eur.nl](mailto:fabbri@law.eur.nl). We are grateful to the Alfred P. Sloan foundation for financial support. We thank Maria Bigoni, Marco Casari, Andrea Geraci, Francesco Parisi, Louis Visscher, Roberto Weber and seminar participants at Erasmus University Rotterdam, University of Bologna, University of Trento, Annual conference of the European Association of Law and Economics 2014, Annual conference of the Italian Society of Law and Economics 2014 for helpful suggestions and Stefano Rizzo for valuable research assistance. The usual disclaimer applies.

## 1. Introduction

Evidence from laboratory experiments suggests that a sizeable fraction of individuals are willing to bear a personal cost for punishing a wrongdoer even if they are not directly affected from the consequences of the rule violation (e.g. Fehr and Fischbacher, 2004b; Balafoutas and Nikiforakis, 2012; Chavez and Bicchieri, 2013). This type of punishment has been labelled "third-party punishment". Social scientists have recently devoted considerable attention to third-party punishment as it can explain the evolution and persistence of social norms in modern large organizations characterized by a majority of anonymous, one-shot interactions (Fehr and Gächter, 2002; Gintis, 2000). However, despite substantial progress in identifying the determinants of decentralized third-party punishment (Bernhard et al., 2006; Buckholtz et al., 2008; Coffman, 2011; Hoff et al., 2011; Lewisch et al., 2011; Lieberman and Linke, 2007; Marlowe et al., 2008; Mathew and Boyd, 2011; Shinada et al., 2004), there is only a partial understanding of what are its major determinants yet (Fehr and Fischbacher, 2004a)<sup>1</sup>.

In this paper we focus on the effects of social influence third-party punishment. By social influence we refer to the effect of the endogenous interactions between a third-party's preferences for punishment and the preferences for punishment expressed by other bystanders (Manski, 2000). Scholars report

---

<sup>1</sup>Contributions in applied psychology investigating the determinants of decentralized TPP have also flourished in the last decades. Kurzban et al. (2007) find, in a laboratory experiment, that subjects increase punishment when observers are present, arguing that TPP is influenced by the so called "audience effect". Subsequent works confirm that anonymity has a causal effect on TPP (Piazza and Bering, 2008), suggesting that the third party decision to sanction wrongdoers is influenced by a cost-dependent reputation effect (Nelissen, 2008) and by emotions (Nelissen and Zeelenberg, 2009). Moreover Lotz et al. (2011) suggest that differences in the level of third-parties punishment provided within a group of agents could be explained by heterogeneity in bystanders' "justice sensitivity".

field and experimental evidence that social influence is a major determinant of human behavior in a variety of settings characterized by important economic consequences<sup>2</sup>. Nonetheless, there is scanty empirical evidence on the effects of social influence on individuals' willingness to engage in third-party punishment. Therefore, it seems relevant to raise the question of whether social influence could affect individuals' willingness to engage in third-party punishment. If that is the case, then policymakers might adjust the level of third-party punishment in a society by making use of social influence effects<sup>3</sup>. This paper aims at filling this gap in the literature on third-party punishment. Particularly, the goal of our paper is twofold. First of all, as stated above, we want to analyze how effective is social influence on the decision of third parties to punish behavior they somehow find wrong or disgraceful. Secondly, we identify the different channels through which social influence could affect third-party punishment. In general, social influence could affect behavior either through the desire that third parties have to conform to the average

---

<sup>2</sup>Starting from the pioneering work of Asch (1951, 1956), contributions in experimental psychology show how individuals tend to modify and distort self-judgments under the influence of group pressure, culture influence and taste for conformism (for a survey see Bond and Smith, 1996). Economists have been mostly interested in the implications of social influence effects for the functioning mechanisms of financial markets. Indeed, most of the contributions focus on the process of information acquisition in investment strategies (Cooper and Rege, 2008; Devenow and Welch, 1996; Scharfstein and Stein, 1990; for a survey see Hirshleifer and Hong Teoh, 2003). Also, scholars in economics investigated the effects of social influence on the labor market. Studies report that peer pressure influences labor productivity (Falk and Ichino, 2006; Mas and Moretti, 2009) and that social networks characterized by an elevated percentage of unemployed individuals could generate social norms perpetuating unemployment (Akerlof, 1980; Topa, 2001). Moreover, reporting results of laboratory experiments, Falk et al. (2010) and Krupka and Weber (2009) find that social influence plays a role in determining pro-social behavior. Empirical evidence suggests that social influence significantly affects teenage pregnancy (Akerlof et al., 1996), obesity (Christakis and Fowler, 2007), judicial voting patterns (Sunstein et al., 2006), investment strategies (Hirshleifer and Hong Teoh, 2003), tax evasion (Fortin et al., 2007; Galbiati and Zanella, 2012) and other criminal activities (Glaeser et al., 1996; Falk and Fischbacher, 2002).

<sup>3</sup>Examples of policies that exploit social influence effects can be found in Cialdini (1993), Coleman (1996) and Perkins et al. (2010).

behavior of other third parties or through the fear of some sort of social sanction in case their punishment choices are disliked by peers. We try to disentangle these two channels, to see whether both affect the decision to punish and whether one is prominent.

We present a laboratory experiment investigating how social influence affects individuals' decision to engage in third-party punishment. The context for our study is a variant of the Fehr and Fischbacher (2004b) dictator game with third-party punishment. Following Cox et al. (2007) and Swope et al. (2008), in our experiment a dictator has the possibility to take from a passive receiver some or all of the initial endowment she received from the experimenter. After observing the dictator's action, a third-party<sup>4</sup> has the opportunity to engage in costly punishment. The game is repeated for two periods and each period consists of two stages. In the first stage of the first period, players are given identical endowments and the dictator has to decide how much to take from the passive receiver he has been matched with. In the second stage, a third-party could reduce dictator's earnings by assigning costly punishment points. After that, the second period starts. Between the two periods, third-party punishers receive information. In this respect, we provide some experimental subjects with information regarding peers' punishment choices in the first period and with information about whether their peers liked or disliked their punishing behavior (*Informational* and *Normative* treatments). Particularly, in what we label *Informational* treatment we provide only information about the average behavior of other third parties in the previous period, whereas in the *Normative* treatment we provide subjects both with information about average behavior and with information about how their decisions were received by peers. The *Informational* treatment is meant to capture the first channel of social influence, that is the desire to conform, the

---

<sup>4</sup>In order to minimize repetitions, when talking about third-parties engaging in punishment we will employ the terms "third-party" and "bystander" interchangeably throughout the paper.

*Normative* treatment tries to capture the extra incentive provided by disapproval. In order to control for possible confounding factors, we also have a *Control* treatment, where only socially irrelevant information is provided to subjects. To identify how social influence affects third-parties' punishment choices we compare behavior when between first and second period they receive socially irrelevant information (*Control*), to behavior in two treatments where they receive information about other bystanders' average punishment. We then compare behavior in the *Information* and *Normative* treatments, to check whether subjects are sensitive to approval and disapproval by peers.<sup>5</sup> Disentangling *Informational* and *Normative* social influence is important for policy purposes because previous studies have stated that, while the former has persistent effects on individual behavior, the impact of the latter is weaker and limited in time (Cialdini and Goldstein, 2004).

Our experiment shows that social influence significantly affects bystanders' decision to engage in third-party punishment. Moreover, we find that third-parties engaging in high level of punishment are the most affected by social influence. We also find that the larger is the difference between a bystander's beliefs regarding peers' punishment and actual peers' punishment, the stronger the bystander's reaction to information about peers' behavior. Finally, we find that subjects may not respond to information in the *Informational* treatment, especially when they disregard their initial beliefs about other peers' behavior in their punishment choices. However, the same subjects tend to reduce the amount of their punishment when they punish more than average and are subject to peers' judgment in the *Normative*

---

<sup>5</sup>The design of our experiment is somewhat similar to Cason and Mui (1998). In their contribution, the authors provide information about the behavior of other peers to experimental subjects. They find that subjects exposed to such information do not change their behavior on average, whereas subjects in a control group who received only irrelevant information become significantly more self-regarding. Furthermore, Deutsch and Gerard (1955) suggest that social influence affects the behavior of agents. Agents derive utility from doing "the right action" and (dis-) utility from being (dis-) liked by her peers.

treatment. Hence, both channels of social influence are important and effective.

The paper is organized as follows. In Section 2 we present the experimental design. Section 3 specifies the theoretical framework and the hypothesis we test. Section 4 presents the experiment results and section 5 discusses the implications of our findings, suggests possible directions for future research and concludes.

## 2. The Experiment

### *2.1. Experimental Design*

The experiment is composed by three treatments. The design of the treatments consists in a variant of the dictator game with third-party punishment, where the dictator has the possibility to take tokens from a passive receiver. The game has 3 possible roles: receiver (Participant A), dictator (Participant B) and Third-party (Participant C). The game is divided in two stages. Each participant is endowed with 30 tokens by the experimenter. In the first stage, Participant B has the possibility to take from 0 up to 30 tokens (in multiples of 5) from A. Participants A cannot undertake any action during the game. In the second stage, Participant C has the opportunity to impose a costly punishment to B. Specifically, C could use up to 20 units of her initial endowment to reduce B's payoff. For each token used by C, the payoff of Participant B is reduced by 4 tokens. Participants C specify how many tokens they use in order to reduce B's payoff for each possible action chosen by B (strategy method). The tokens C uses for punishment and the consequent reduction of B's payoff have no effect on the payoff of player A. Agents have full information regarding the rules of the game.

Before the game starts, participants' beliefs about the average punishment choices of the peers are elicited. To do so, we use an incentivized coordination game similar to Krupka and Weber (2013). We refer to this part of the experiment as the "Beliefs elicitation game". We present to participants

an hypothetical situation identical to the game described above. We ask each participant to indicate, for each of the 7 possible actions of B, the number of tokens [0; 20] that in their opinion C would use to punish B. We explain that, once each participant present in the laboratory has provided her answers, the computer selects one of the seven possible actions of B. For the selected action, a participant earns 40 tokens if the number she indicated is equal, bigger or smaller by one unit to the average number indicated by all the participants in the experimental session. Therefore, in this part of the experiment participants have incentives to reveal their true beliefs regarding peers' choices of punishment.

We propose 2 effect treatments (*Informational* and *Normative*) and a control treatment (*Control*). The experiment is composed by two periods. The elicitation of beliefs and the first period of the game are identical in all treatments. Specifically, the amount of tokens B decides to take from A in the first stage is not immediately observed by C. Instead in stage 2 of the first period C's decisions of punishment are elicited employing the "strategy method": for each possible action of B, C states his decision of punishment. Participants are informed that only the punishment decision corresponding to the actual choice made by B determines payoffs. The punishment tokens used by C in correspondence of the other possible choices of B do not have payoff consequences. First period earnings and choices of peers are not revealed to participants at the end of the first period.

At the beginning of the second period, participants' endowments are restored at the initial level. Earnings of the first period are independent from those of the second. The first stage of the second period is identical to the first stage of the previous one: B is endowed with the same amount of tokens and may take part or all of A's endowment. Also in the second stage of the second period, C has to indicate the level of punishment inflicted for each possible action of B. The difference between treatments consists in the kind and amount of information disclosed to participants C before the punishment choices.

In the *Informational* treatment, each participant C receives information about the average number of tokens used to punish B in the first period by the participants C taking part to the experimental session.

In the *Normative* treatment, each participant C receives the same information of *Informational*. However, she is additionally informed that her punishment decisions of the second period will be revealed to 5 peers randomly selected among the experiment participants. After observing these choices, the 5 peers vote for sending an emoticon that will appear to the screen of the participant C. The 5 peers could vote for a smiling emoticon or a sad emoticon. If the majority vote for a smiling emoticon, on player C's monitor will appear a smiling emoticon. A sad emoticon will appear on the screen otherwise. Participants are informed that the emoticon received has no effect on earnings and that it disappears after one minute.

In the *Control* treatment, no relevant information about participants punishment choices is disclosed in the second period. However, we have to rule out the possibility that a change in punishment behavior between periods is driven by factors other than the exposure to social influence. One possible confounding factor is subjects' experience that increases between periods. Another possible confounding factor is that processing new information imposes a cognitive effort to subjects. Hence these factors, not the social content of the information received, could be responsible for a modification of the punishment choice. In fact, in the *Informational* and *Normative* treatments subjects have to process some sort of information, and there is evidence that individuals exposed to a cognitive load tend to modify their behavior (for discussion on this point see Cason and Mui, 1998).

In order to rule out these confounding factors and isolate social influence effects, we then expose the *Control* group to some social irrelevant information. Specifically, we ask at the beginning of the session to each participant her day of birth  $\in [1; 31]$  and we take the average. Since we do not ask nor we report them neither the year nor the month of birth, reporting this measure



do not convey any relevant social information. However, in this way participants in *Control* are affected by the same cognitive burden of participants in treated groups and the only difference lies in the exposure to relevant social information.

In all treatments, the per period payoffs are calculated as:

- $\Pi_A = 30 - t$
- $\Pi_B = 30 + t - 4*p$
- $\Pi_C = 30 - p$

where:

- $t$  = tokens taken by B from A
- $p$  = punishment tokens used by player C

In order to calculate individual earnings, participants are randomly divided in groups of 3. Each group is composed of one participant A, one B and one participant C. The final payment for each group follows this procedure:

- In the dictator game, after the second period is concluded, one of the two periods is randomly selected. This period will be called the "payment period".
- For each participant, earnings relative to the payment period are added to earning collected in the beliefs elicitation game. Earnings from the non-selected period are not paid out.
- A 5 euro participation fee is added to total payments.

*The Procedure.* The experiment was programmed using the software z-Tree (Fischbacher, 2007). Every session was conducted at the Bologna Laboratory for Experimental Social Sciences at the University of Bologna, Italy, between

November and December 2013. All sessions were run by the same experimenter. Participants were for the vast majority graduate and undergraduate students of the University of Bologna, plus some private citizens, and were recruited through the online system ORSEE (Greiner, 2004). In each session participants were split into 5 groups of 3 subjects<sup>6</sup>. Overall, 9 sessions were run, 3 for each treatment, that results in a total of 132 participants (56% female).

In each session, before each of the three part of the experiment (elicitation of beliefs, first period punishment and second period punishment), a printed copy of the instructions was distributed and read aloud by the experimenter<sup>7</sup>. Participants had additional images and tables summarizing the instructions on their computer screen. Information regarding payoff functions and rules of the game was common knowledge. Participants had the possibility to ask questions before the experiment started.

At the end of each session participants completed a brief socio-demographic questionnaire<sup>8</sup>. Each participant took part in one session only. Peers' identities were maintained unknown even after the end of the experiment. In order to guarantee anonymity, participants were individually and privately paid after the experiment finished. No communication among participants was allowed.

The part of the session concerning beliefs elicitation and treatments lasted around 20 minutes. However, due to the impossibility of learning throughout periods and the limited number of decisions each participant had to take, we were concerned about the possibility that instructions may were not fully

---

<sup>6</sup>In one session of the *Informational* treatment there were only 4 groups, for a total of 12 participants.

<sup>7</sup>Original instructions are in Italian and are available upon request. A copy of the instructions for the *Normative* treatment translated in English is included in Appendix B.

<sup>8</sup>In one session of the *Informational* treatment subjects' socio-demographic characteristics were not recorded due to a technical problem.

understood. In order to minimize this possibility, we adopt special care in writing detailed instructions and providing multiple examples, and we also asked subjects to correctly answer control questions before proceeding with each part of the experiment. As a result, each experimental session lasted in total about 45 minutes. Tokens were converted into euro at a rate of 5 tokens for 1 euro. Subjects earned on average approximately 11 euros for experimental session.

### 3. Hypothesis

Following the customary assumptions, the predictions of the game outcomes are straightforward. Agents' utility is an increasing function of individual wealth and agents are individual payoff maximizer. Hence, in any treatment, no punishment should be observed, since the payoff-maximizing strategy for third-parties is to punish nothing and keep the initial endowment. Anticipating the absence of punishment, dictators should take all the tokens from receivers.

However past dictator game experiments have shown two behavioral regularities. On one hand, even in games where the dictator faces no threat of punishment, positive amounts of tokens are transferred (in our setting: are left) to the receiver. On the other hand, third-parties engage in costly punishment for dictator's levels of transfer (in our setting: for dictator's levels of taking) perceived as unfair. In this study we are interested in verifying how, given the action of a dictator, the punishment choices of other third-parties affects the utility that a bystander derives from punishing the dictator.

Consider the choice of a third-party  $i$  to use  $p$  tokens of her initial endowment in order to punish a dictator that takes  $z$  tokens from a passive receiver. Third-parties' individual utility is an increasing function of the final monetary earnings  $x$ . Moreover, given a dictator's action, third-parties have some inherent preferences  $p_z^k$  for the amount of tokens she wants to use for punishment.  $p_{i,z}^k$  could be interpreted as reflecting the individual sense of justice of

the third-party *we*. If a third-party chooses to punish the dictator a quantity different from her inherent preference, she has to bear a cost  $s$  that increases when the absolute difference between  $p^k$  and the  $p$  increases.

Furthermore, third-parties have some beliefs  $E(\bar{p})$  regarding the average amount of tokens that the other bystanders will use for punishing dictators. A third-party incurs a cost  $c$  for punishing a quantity of tokens different from  $E(\bar{p})$ , and this cost becomes larger when the absolute difference between individual punishment and the average punishment of the peers increases.  $c$  incorporates both the costs imposed by the other bystanders observing the third-party that deviates from the average punishment and the disutility the third-party experiences in not conforming with the peers' behavior independently from the fact that her action is observed.

Therefore, in their punishment decisions a third-party maximizes individual utility taking into account the cost of using tokens for punishing a dictator and so reducing her monetary payoff, the cost for deviating from her inherent preference for punishment and the cost of not conforming to the peers' average punishment:

$$\begin{aligned} \max_{p_{i,z,t}} \quad & U_{i,z,t} = x_{i,z,t} - (s(E_i(\bar{p}_{z,t}) - p_{i,z,t})^2 + c(p_{i,z}^k - p_{i,z,t})^2) \\ \text{s. t.} \quad & y_w e = x_{i,z,t} + p_{i,z,t} \end{aligned} \tag{1}$$

Where  $y$  is third-party's initial endowment. Assuming an interior solution exists, equation (1) generates the following first order conditions:

$$p_{i,z,t}^* = \frac{sE_i(\bar{p}_{z,t}) + cp_{i,z}^k - 1}{s + c} \tag{2}$$

Therefore, according to our model of social influence the optimal punishment choice of a third-party is an increasing function of the expected punishment chosen by her peers. Furthermore, the higher is for a third-party the cost  $s$  of not conforming to other bystanders' average punishment relative to the cost  $c$  of deviating from inherent preferences, the higher will be the tendency to

conform to the peers' average punishment. Allowing for concavity of agents' utility function, the intuition of the previous results will still work.

In order to test our predictions, as a first step we verify if there is a positive association between a third-party's beliefs regarding peers' average punishment and her first period punishment. As a second step, we then investigate how participants modify their punishment choices between first and second period. Assume that third-parties in TREATED revise beliefs about average peers' punishment substituting their initial priors with the actual punishment level revealed them after the first period, hence  $E_i(\bar{p}_{z,t}) = \bar{p}_{z,t-1}$ . The punishment variation across periods is given by:

$$(p_{i,z,2}^* - p_{i,z,1}^*) = \frac{s(\bar{p}_{z,1} - E_i(\bar{p}_{z,1}))}{s + c} \quad (3)$$

For the moment, focus only on the distinction between participants in *Control* and the other participants grouped together, that we call TREATED. Third-parties in *Control* are not exposed to social relevant information between period 1 and 2. Instead, bystanders in TREATED are exposed to information that may induce them to update their initial beliefs regarding peers' average punishment and so influence their second period punishment decision. Therefore, if social influence has an effect on third-party punishment decision, we expect participants in TREATED to be more likely to modify their punishment decisions in the passage between the first and second period punishment compared to participants in *Control*.

Thus, according to our model revealing to a bystander her peers' average punishment may trigger a change in her second punishment decision as a consequence of a beliefs updating process. Specifically, for a bystander the likelihood to change punishment decision in the second period increases when the absolute difference between her beliefs regarding peers' average punishment and the actual average punishment of the first period is large. Therefore, we test the following hypothesis:

1. Zero Social Influence hypothesis: *In the first period, punishment decisions of bystanders are not influenced by their beliefs regarding peers' average punishment. Moreover, bystanders in TREATED are as likely as bystanders in Control to modify their initial punishment decisions.*

Second, we want to identify who are the bystanders more responsive to social influence. Third-parties deciding to use tokens for punishing a dictator are reducing their final monetary payments. Hence, every time we observe a bystander punishing a positive amount, according to our model we infer that  $sE(\bar{p}) + cp^k - 1$  is positive. This could mean that the bystander has inherent preferences for punishing a positive amount ( $p^k > 0$ ) and at the same time she attaches a positive weight to this component of the utility function ( $c > 0$ ). However, it is also possible that the bystander attaches a positive weight to the social component of the utility function ( $s > 0$ ) and she expects peers to punish on average a positive amount of tokens ( $E(\bar{p}_{z,t}) > 0$ )<sup>9</sup>. If this last possibility is true, the higher is a bystander's punishment in the first period the more she attaches weight to the social component of the utility function and so the more likely she is to modify second period punishment decision. Now consider the difference between first and second period punishment of a bystander. Inherent preferences for punishment are stable, so they do not play a role in the decision to eventually modify punishment choice. Instead, according to the prediction of our model, the larger is  $s$  for a bystander, the more she responds to the social information regarding peers' punishment. Therefore, holding constant  $E(\bar{p}_{z,t}) - \bar{p}_{z,t-1}$ , we expect that the more a bystander punished in the first period, the more she is likely to revise her punishment decisions in the second period.

Moreover, we also consider the difference between a bystander's beliefs re-

---

<sup>9</sup>Of course, it is possible that what it is observed is a combination of this two possibilities.

garding peers' average punishment and her first period punishment. In the first period, a bystander could punish an amount different from her beliefs regarding peers' average punishment because she only cares about her monetary payoff or because her inherent preference for punishment differs from the expected average punishment and the cost  $s$  of non conforming to peers' average punishment is small compared to the cost  $c$  of non following inherent preferences. In both cases, the choice of the bystander reveals that in her punishment decisions she is little influenced by peers' behavior. As a consequence, we expect the more a bystander punishes in the first period a quantity different from her beliefs about peers' average punishment, the less she will be responsive to social influence.

2. Differential Social Influence hypothesis: *Third-parties that engage in high punishment in the first period are the most responsive to social influence. Conversely, the higher the absolute difference between a bystanders' beliefs regarding average peers' punishment and first period individual punishment, the lower the bystander likelihood to modify punishment choices in the second period.*

Finally, we investigate the psychological mechanisms triggering social influence. In our experiment, we give bystanders in *Normative* and *Informational* the same information about peers' punishment. However, in the *Informational* treatment, the second period choices of the third-party are not observable ex post by other participants. As a consequence, in the *Informational* treatment a third-party has no incentives to conform to peers' punishment choices if her only goal is being liked by them. Hence, a bystander would eventually modify his punishment strategy only if *Informational* social influence is at work.

On the other hand, in the *Normative* treatment a bystander is aware that her punishment decisions of the second period will be observed by peers and that

they will express a judgement regarding those choices. As a consequence, for bystanders in *Normative* the cost  $s$  of not conforming to the average peers' punishment has been modified in the passage between first and second treatment since we added a *Normative* social influence component. Hence, if some bystanders are responsive to *Normative* but not to *Informational* influence, the *Normative* treatment will show social influence effects different from those resulting from the *Informational* treatment. The difference in the way bystanders modify their punishment decisions between *Normative* and *Informational* treatments isolates the effect of *Normative* social influence on third-party punishment.

3. Equivalence of *Normative* and *Informational* Influence hypothesis: *Social influence effects on third-party punishment are the same for subjects exposed to Informational and Normative influence.*

#### 4. Results

Table 1 reports summary statistic relative to our data<sup>10</sup>. Dictators leave approximatively 36% of receivers' endowment. This finding is consistent with results from other comparable experiments where dictators have to take tokens from the endowment of a passive receiver (List, 2007; Krupka and Weber, 2013)<sup>11</sup>.

---

<sup>10</sup>Additional summary statistics where we consider separately the 7 possible punishment choices in each period, are reported in Tables A.7 and A.8 in the appendix.

<sup>11</sup>In the classical dictator game without punishment a dictator has the possibility to give part of his endowment to a passive receiver. In a meta-study, Engel (2011) found that on average dictators give roughly 25% of their endowment to the receiver. However, in our design dictators has to *take* money from receivers' endowment instead of *giving* them. This difference and the possibility of being punished that characterizes our design are likely to explain the slightly more fair allocation we registered compared to the standard dictator game (on this point, see also Krupka and Weber, 2013).



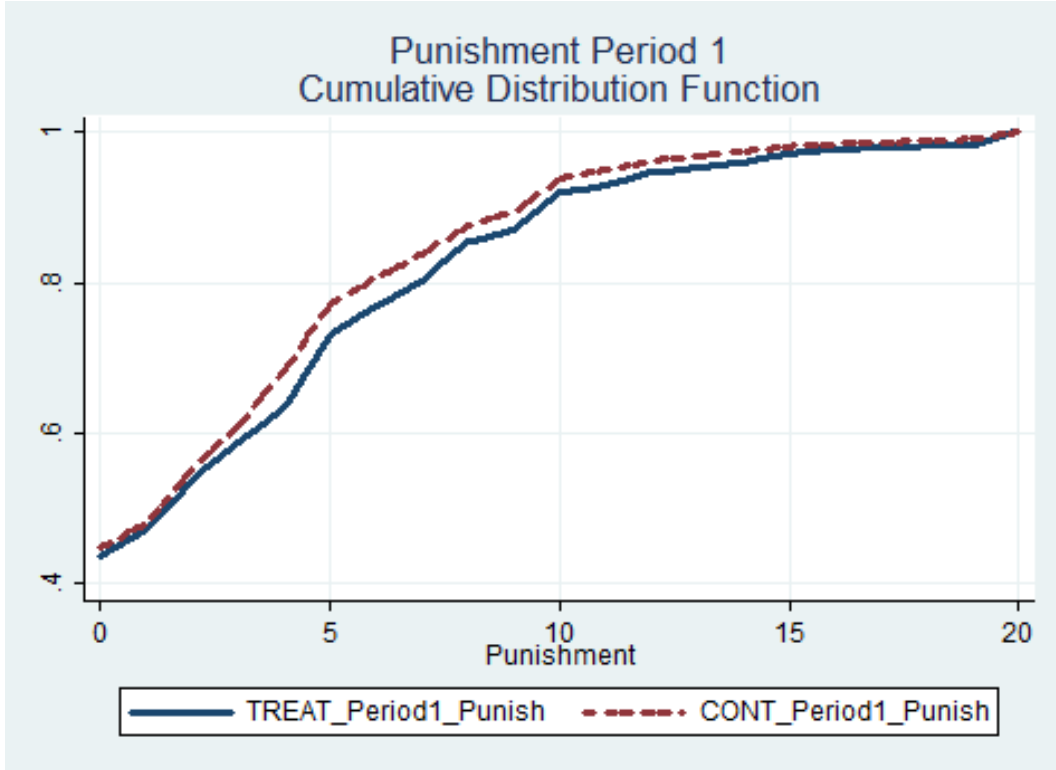


Figure 1: Punishment Period 1 Cumulative Distribution Function

On average bystanders punish approximately 3.5 tokens, decreasing punishment amount in the second period. When the dictator takes all the money from the receiver, third parties spend approximately 6 tokens in punishment. Average punishment then progressively decline, reaching virtually 0, when levels of dictators' taking decrease. Also this result is consistent with previous findings on third-party punishment in dictator games (Fehr and Fischbacher, 2004b). However, if we consider the 3 treatments separately, we see that in *Control* punishment slightly increases in the second period, while in both *Normative* and *Informational* it decreases. Considering third-parties' beliefs regarding peers' punishment behavior, beliefs are on average are higher than actual punishment.

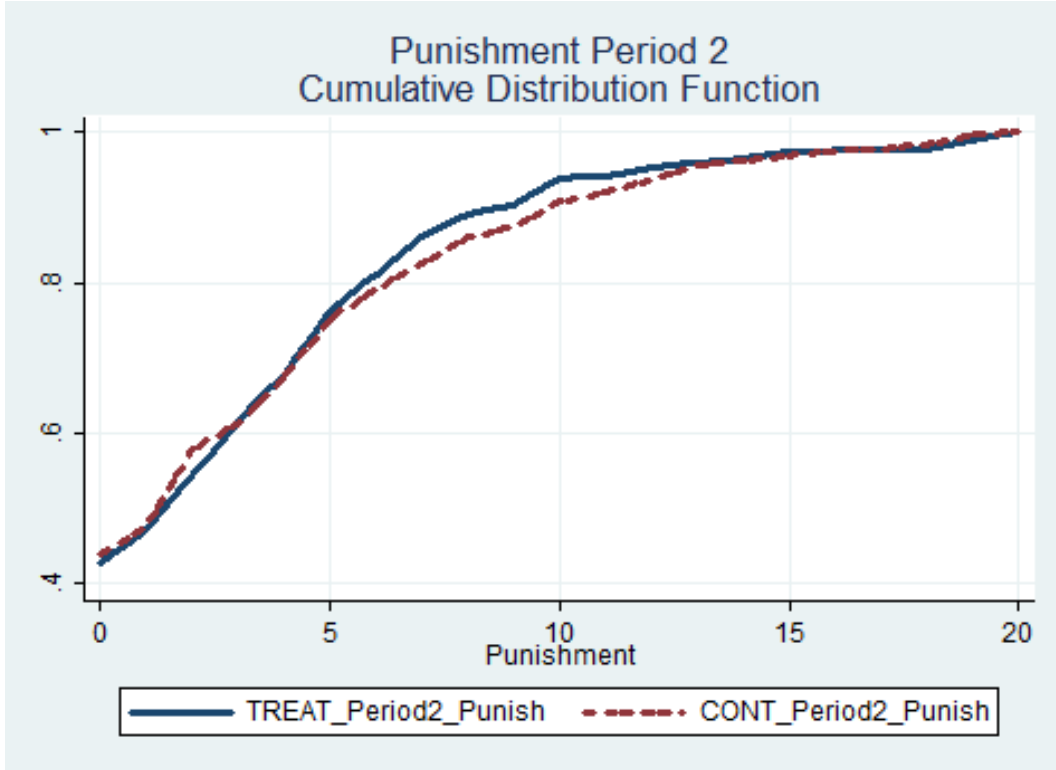


Figure 2: Punishment Period 2 Cumulative Distribution Function

we proceed considering, for each bystander in a single punishment period, the average of her 7 punishment choices corresponding to different levels of dictator taking. we compare the cumulative distribution of this measure in TREATED and *Control*. Figures 1 and 2 report the cumulative distribution functions of punishment respectively in period 1 and 2. In the first period, the cumulative punishment choice distribution in *Control* exceeds the distribution in TREATED for any possible punishment level. However, a Kolmogorov-Smirnov test cannot reject the null hypothesis that the distributions are equivalent. In the second period instead, the cumulative punishment choice distribution of TREATMENT exceeds the distribution of *Control* for some punishment levels greater than 5. However, also in the

second period a Kolmogorov-Smirnov test cannot reject the equivalence of the two distributions, nor the samples means are statistically different (t-test two sided p-value 66%). Now we test our hypothesis.

#### 4.1. Zero Social Influence Hypothesis

We start by investigating if bystanders punish in the first period according to their beliefs regarding the average punishment they expect peers' will use. As a first step, we test if there is a significant difference between the punishment used by a bystander and her beliefs about peers' average punishment. We conduct a t-test comparing the two averages under the null hypothesis that they are the same. Third-parties punish on average 3.6 tokens in the first period while their beliefs about peers' average punishment is 5.1 tokens.

Results of the t-test reject our hypothesis and indicate that bystanders in the first period punish significantly less than what they think peers on average will do (t-test two-tails, p-value < 1%).

We want to verify if this result is driven by those third-parties that during the experiment punish always 0 (we name them "selfish"). Excluding selfish punishers from the sample, bystanders' beliefs about peers' punishment remain higher than the punishment they provide (6.0 versus 5.4 tokens), however the difference is not statistically significant (p-value 12%). This results suggest that selfish subjects are responsible for the aforementioned gap.

we continue the analysis regressing the quantity of punishment tokens a bystander uses in the first period with her beliefs regarding peers' average punishment and a set of socio-demographic characteristics. Results are reported in Table 2<sup>12</sup>. The variable *Beliefs\_Punish* indicates bystanders' beliefs about peers average punishment. The coefficient is positive and significant at the 1% level in any model specifications. According to the model estimation, bystanders spend additional 0.4 token for every unit of increase in expected average peers' punishment. Hence, data suggest that third-parties

---

<sup>12</sup>Table C.9 in Appendix C reports a description of each variable employed.

Table 1: Summary Statistics

<b>Treatment</b>	<b>male</b>	<b>age</b>	<b>dictatorTake</b>	<b>Beliefs</b>	<b>PunishPer1</b>	<b>PunishPer2</b>
<i>Control</i>						
(Mean)	.33	24.68	18.38	5.08	3.33	3.54
(Median)	0	25	20	4.29	2.71	2.71
(SD)	.48	2.76	11.19	3.71	3.42	3.82
<i>Normative</i>						
	.58	26.18	18.28	4.80	3.30	3.03
	1	25	20	4.57	3	2.29
	.50	5.38	11.40	3.40	3.34	3.37
<b>Informational</b>						
	.37	25.15	17.32	5.41	4.15	3.81
	0	25	16.25	4.86	4.29	3.86
	.49	3.72	10.85	3.67	3.71	3.90
<b>Total</b>						
	.44	25.37	18.01	5.09	3.58	3.45
	0	25	20	4.71	3.43	3.07
	.50	4.17	11.08	3.57	3.48	3.69

are influenced in their first period punishment decisions from beliefs about peers' punishment. This finding goes against the Zero Social Influence hypothesis.

we proceed in the analysis verifying how third-parties in *Control* and in *TREATED* modify punishment choices between first and second period. As a first step, we sort bystanders into 2 main categories: those who never change punishment decisions across periods and those who change at least once. In

Table 2: Determinants First Period Punishment

	(1)	(2)
Beliefs	0.405*** (0.05)	0.432*** (0.06)
male	-0.515 (0.58)	-0.407 (0.62)
age	-0.078 (0.06)	-0.047 (0.07)
degree	-0.842** (0.41)	-1.435*** (0.45)
worker	0.759 (0.74)	0.860 (0.80)
social	0.378 (0.73)	0.399 (0.85)
arts	0.234 (1.06)	0.021 (1.28)
field_other	0.409 (0.63)	-0.131 (0.75)
risk	0.138 (0.10)	0.037 (0.11)
logic	-0.305 (0.40)	0.057 (0.47)
impulsivity	-0.636** (0.27)	-0.687** (0.32)
Instruction	0.000 (0.00)	-0.000 (0.00)
DictatorTake	-0.007 (0.01)	0.026 (0.02)
_cons	6.246*** (1.90)	7.889*** (2.04)
<i>N</i>	924	630
<i>R</i> <sup>2</sup>	0.281	0.275
<i>BIC</i>	5203.7	3602.6

Notes: OLS regression: dep. var. *Strat\_Punish*, SE clustered by subject. Significance levels: \* p<0.10, \*\* p<0.05, \*\*\* p<0.01

the TREATED group 53 subjects (61%) change at least one punishment decision between periods, while in the *Control* treatment 24 subjects (53%) do so. This difference is not statistically significant. If we repeat the same test excluding selfish bystanders, it turns out that in TREATED 87% and in *Control* 80% of third-parties change punishment decisions at least once. However, also in this case the difference is not statistically significant.

We also verify how many times on average each punisher change decision across periods. In TREATED bystanders change decision 2.5 times, while in *Control* they change 2.1 times. This difference is not statistically significant and it remains roughly unchangend even if we exclude selfish bystanders. Therefore, these results do not provide evidence against the Zero Social Influence hypothesis. The result seems to be driven by the high percentage of participants (53%) in the *Control* group that modifies punishment choices at least once, even if they did not receive any relevant social information.

As a second step, we test if there is a difference in the likelihood that participants in *Control* and TREATMENT change punishment decisions. We create the dummy variable *DummyP1p0* that takes the value 1 when punishment in the second period differs from punish in the first one and 0 otherwise. We implement a logistic regression to estimate the likelihood of changing punishment choice across periods. Results of the model are presented in Table 3. The dummy *TREATED* equal to 1 for participants in *Normative* and INFORMATIONAL. The coefficient of the dummy is positive and statistically significant in any of the model specifications<sup>13</sup>. Therefore, we conclude that the results of the logistic regression do not support the Zero Social Influence hypothesis and indicate that participants in TREATED modify punishment decisions across periods more often than those in *Control*.

---

<sup>13</sup>Model 2 differs from Model 1 because excludes selfish participants. Model 3 add the *Control* variables *Strat\_Punish*, indicating punishment exerted in the first period, and *Absp0Belifs*, reporting the absolute difference between a bystander’s punishment in the first period and her beliefs regarding peers’ average punishment. Model 4 excludes selfish participants from the sample.

As a third step, we investigated how third-parties modify their punishment choices. In *Control* bystanders reduce punishment in the second period 48 times (15%), increase 49 times (16%) and do not change 218 times (69%). In TREATED, bystanders reduce punishment 140 times (23%), increase 93 (15%) and do not change 376 times (62%). This choices result for *Control* in an average increase in punishment from period 1 to period 2 of 0.21 tokens (from 3.3 to 3.5) and in an average decrease in TREATED of 0.30 tokens (from 3.7 to 3.4). The mean punishment difference across periods is not statistically different in *Control* and TREATED (p-value 0.16, t-test two-sided).

However, we could expect that bystanders have no reason to punish a dictator when she does not take any amount of money from the receiver. Hence, when the dictator chooses to take 0 from the receiver, we expect little or no punishment both in *Control* and TREATED. In fact, if we exclude the situations in which the dictator takes 0 from the receiver, the average difference between bystanders' punishment in the two periods is weakly statistically significant higher in *Control* versus TREATED (p-value 0.09). Furthermore, if we consider only situations in which the dicator takes half or more of receiver's initial endowment, this difference between *Control* and TREATMENT becomes significant at the 5% level.

Hence, results of this third set of tests suggest that, at least for situations where dictators subtract positive amounts of tokens from receivers, bystanders exposed to social influence significantly reduce the amount of punishment provided compared to bystanders in *Control*. These results provide evidence against the Zero Social Influence hypothesis.

Fourth, we test the hypothesis that a large absolute difference between a bystander's beliefs regarding peers' average punishment and actual peers' average punishment increases the likelihood to modify the initial bystander punishment choice. Third-parties receive information regarding actual peers' punishment in the *Normative* and *Informational* treatments only, so we

restrict the analysis to these treatments. we test this hypothesis using a logistic model. Results are reported in Table 3. From models 7 and 8, we can see that the coefficient of the variable *Abs\_BeliefAvgPunish* is positive as expected, however only weakly significant. The estimations suggest that an increase of one unit in the difference *Abs\_BeliefAvgPunish* increases on average the probability of modifying second period punishment by 3.5%<sup>14</sup>. This result provides evidence against the Zero Social Influence hypothesis. Finally, we report a series of statistical tests that account for the direction of punishment deviation between periods. First, we consider the variable *P1p0* that reports the difference between punishment in period 2 and period 1. Table 4.1 provide summary statistics of this variable, grouping subjects according to the difference between individual beliefs about average punishment and actual average punishment observed. When beliefs match exactly the average punishment observed ( $P1p0\_BAP = 0$ ), subjects confirm punishment choices in the second period 76% of times. Instead, when beliefs are larger or smaller than the actual average punishment observed ( $P1p0\_BAP > \text{ and } < 0$  respectively), subjects confirm first period choice respectively 51% and 73% of times. Considering how subjects modify their decisions, we see that those observing actual average punishment smaller than beliefs reduce on average punishment in the second period of 0.30 tokens. Instead, subjects observing actual average punishment equal to beliefs reduce punishment of 0.14 tokens between the two periods, and those observing average punishment smaller than beliefs increase punishment in the second period of 0.27 tokens.

---

<sup>14</sup>we also consider the possibility that a large absolute difference between a bystander punishment in the first period and the average peers' punishment increases the likelihood to change punishment decision in the second period. We create the variable *Abs\_Signalp0*, reporting the absolute difference between individual punishment in the first period and average peers' punishment. Results of the logistic estimations are reported in model 7 and 8 of Table 3. The coefficient of *Abs\_Signalp0* is not statistically different from 0 in any model specification. As a consequence, we conclude that *Abs\_Signalp0* has no impact on subjects' likelihood to modify punishment decision.



Table 3: Probability modify punishment across periods

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
male	-0.219*** (0.08)	-0.258*** (0.09)	-0.168** (0.07)	-0.215** (0.09)	-0.098 (0.12)	-0.019 (0.14)	-0.089 (0.12)	-0.017 (0.14)
risk	0.030* (0.02)	0.014 (0.02)	0.019 (0.02)	0.008 (0.02)	-0.015 (0.02)	-0.040 (0.03)	-0.020 (0.02)	-0.043* (0.03)
logic	-0.077 (0.06)	-0.033 (0.07)	-0.061 (0.06)	-0.031 (0.07)	0.135 (0.09)	0.242** (0.10)	0.149 (0.09)	0.246** (0.10)
TREATED	0.155** (0.06)	0.131* (0.07)	0.124** (0.06)	0.122* (0.07)				
Strat_Punish			0.040*** (0.01)	0.021** (0.01)	0.044*** (0.01)	0.025** (0.01)	0.047*** (0.01)	0.029** (0.01)
Abs_p0Belifs			-0.010 (0.01)	-0.003 (0.01)	-0.026** (0.01)	-0.019* (0.01)	-0.032** (0.01)	-0.026 (0.02)
Abs_Signalp0							0.002 (0.01)	0.006 (0.01)
Abs_BelAvgPun							0.024 (0.02)	0.023 (0.03)
Beliefs			0.015 (0.01)	0.009 (0.01)	0.025** (0.01)	0.021* (0.01)	0.015 (0.01)	0.009 (0.01)
Other contr	Y	Y	Y	Y	Y	Y	Y	Y
<i>N</i>	924	630	924	630	609	420	609	420
pseudo <i>R</i> <sup>2</sup>	0.092	0.059	0.207	0.095	0.206	0.101	0.226	0.123

Notes: Logistic regression: dep. var. *DummyP1p0*, marginal effect at means, SE clustered by subject. Significance levels: \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. Other *Controls* include: degree worker social arts field\_other DictatorTake impulsivity age.

Additionally, we want to test if subjects that receive information regarding session average punishment modify their second period punishment decisions in order to conform to peers' average punishment. We consider each participant's punishment decision in the first and in the second period (*P1* and *P2* respectively) and the information regarding session average first period punishment that the TREATED group receive between first and second period (*AvgPunish*). If a subject wants to conform her punishment behavior to the punishment choices of the peers, she will choose a second period pun-

ishment  $P2$  that reduces the gap between first period punishment  $P1$  and the revealed peers' average punishment  $AvgPunish$ . Therefore, we create the variable  $Diff\_Deviation$ , that compares for each punishment choice the absolute value of the difference between participants' first period average punishment and individual punishment in the second and in the first period:  $Diff\_Deviation = |P2 - AvgPunish| - |P1 - AvgPunish|$ . In *Control*, where subjects are not exposed to information regarding peers' punishment behavior,  $Diff\_Deviation$  increases of 0.20 unit. Conversely,  $Diff\_Deviation$  is on average reduced of 0.25 unit for subjects belonging to the TREATED group. A t-test test comparing the mean values of  $Diff\_Deviation$  between groups under the alternative hypothesis that this measure is smaller in TREATED than in *Control* reports weakly statistical evidence against the null hypothesis (t-test one side, p-value 0.06).

Table 4: Punishment difference across periods

<b>Treatment</b>	<b>P1p0</b>	<b>P1p0_BAP&gt;0</b>	<b>P1p0_BAP=0</b>	<b>P1p0_BAP&lt;0</b>
<b>Informational</b>	294	142	27	125
(mean)	-.33	-1.04	-.11	.42
(sd)	3.28	3.96	.42	2.52
<i>Normative</i>	315	170	60	85
	-.28	-.49	-.15	.06
	2.58	3.09	1.63	1.90
<b>Total</b>	609	312	87	210
	-.30	-.74	-.14	.27
	2.94	3.52	1.37	2.29

Notes: Variable  $P1p0$  indicate the difference between punishment in first and second period.  $P1p0\_BAP >, <, = 0$  indicate respectively  $P1p0$  when individual beliefs regarding average punishment are  $>, <, =$  actual average punishment.

Finally, we test the hypothesis that subjects conform to peers' average punishment by means of a logistic regression. First, we create the dummy variable  $ConvergeDummy$  that takes value 1 when  $Diff\_Deviation$  is positive,

that is when a subject chooses a second period punishment that reduces the gap between her first period punishment and peers' average punishment. Then we implement a logistic regression to test if the probability of reducing this gap is the same for participants in *Control* and TREATED groups. Results are reported in Table 5. As we can see from model (1), the coefficient of the dummy TREATED is positive and statistical significant at the 1% level. This result suggests that, if compared to participants in the *Control* group, subjects exposed to information regarding peers' average punishment are roughly 14% more likely to modify their second period punishment choice in order to conform to peers' punishment behavior.

- Conclusion relative to the Zero Social Influence hypothesis: *Bystanders' punishment choices in the first period are positively associated with their own beliefs about average peers' punishment. However, we find that on average bystanders punish less than the expected average peers' punishment. This result seems to be driven by those bystanders that always decide not to punish dictators. We also found evidence that subjects in TREATED are more likely to change punishment decision across periods. Moreover, in Control the amount of punishment in the two periods remains constant, while in TREATED it decreases. The mean punishment difference across periods is statistically higher in Control than in treated if we consider situations where dictators take positive amounts from receivers' endowment. Additionally, a large absolute difference between a bystander's beliefs regarding peers' punishment and actual peers' average punishment increases her likelihood to modify punishment decisions across periods. Finally, if compared to participants in Control, subjects in TREATED choose more often a second period punishment that reduces the gap between individual first period punishment and peers' average punishment.*

*These results provide evidence against the Zero Social Influence hypothesis. Hence, we conclude that the Zero Social Influence hypothesis is*

Table 5: Probability Second Period Individual Punishment Converges Toward Average Participants' First Period Punishment

	(1)	(2)	(3)	(4)	(5)
TREATED	0.144*** (0.05)				
<i>Normative</i>		0.021 (0.08)	-0.018 (0.08)	0.201* (0.11)	0.226** (0.10)
male	-0.089* (0.05)	-0.082 (0.09)	-0.089 (0.08)	-0.195 (0.12)	-0.224** (0.11)
Strat_Punish	0.047*** (0.01)	0.046*** (0.01)	0.049*** (0.01)	0.020*** (0.01)	0.041*** (0.01)
Beliefs	0.000 (0.01)	0.004 (0.01)	0.002 (0.01)		0.010 (0.01)
Punish	-0.033*** (0.01)	-0.023* (0.01)	-0.025* (0.01)		-0.030** (0.01)
Abs_BeliefAvgPunish	0.003 (0.01)	0.014 (0.02)	0.032 (0.02)		-0.010 (0.01)
Abs_p0Belifs			-0.045*** (0.02)		-0.012 (0.01)
NormAbs_BeliefAvgPunish			-0.025 (0.02)		
NorAbs_p0Belifs			0.044** (0.02)		
Other contr	Y	Y	Y	Y	Y
<i>N</i>	924	609	609	197	197
pseudo $R^2$	0.195	0.163	0.185	0.198	0.310

Notes: Logistic regression: dep. var. *ConvergeDummy*, marginal effect at means, SE clustered by subject. Significance levels: \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. Other *Controls* include: age risk impulsivity logic student.

*not supported by the results of the experiment.*

#### 4.2. Differential Social Influence hypothesis

First, we want to verify if bystanders that engage in less punishment in the first period are also less responsive to social influence. Third-parties receive relevant social information in the *Normative* and *Informational* treatments only, so we restrict the analysis to these treatments. For each

of the 7 dictators' decisions, we characterize bystanders that punish in the first period above the median as "high punishers"<sup>15</sup>. For each transfer level considered, the percentage of third-parties modifying punishment decision across periods among high punishers is almost double than among the other punishers. If we exclude selfish bystanders we still have similar results.

we test the hypothesis implementing a logistic model. We estimate the probability of modifying punishment decision including the independent variable *Strat\_Punish* that reports the level of punishment provided in the first period. Results are reported in Table 3. In any model specification, the coefficient associated with *Strat\_Punish* is positive and significant at the 1% level. The coefficient of *Strat\_Punish* suggests that, holding constant at their means the other *Controls*, a bystander spending 1 additional token in first period punishment is 3% to 5% more likely to revise her punishment choice in the second period.

Hence, we conclude that this first set of results supports the Differential Social Influence hypothesis.

Second, we want to verify if bystanders that choose to punish in the first period a quantity different from their beliefs regarding average peers' punishment are less responsive to social influence compared to the other bystanders. we implement a logistic regression estimating the probability that a bystander modifies punishment decisions across periods. As independent variable, we introduce *Abs\_p0Belifs*, the absolute difference between a bystander's beliefs regarding average peers' punishment and her individual punishment in the first period. We report results in Table 3.

The coefficient of *Abs\_p0Belifs* is negative and statistically different from 0 in model specifications 7 and 8, in which we include all the *Controls*<sup>16</sup>.

---

<sup>15</sup>Results are substantially the same if we choose the average punishment as a criterion for classification.

<sup>16</sup>In model 8 we exclude from the analysis selfish bystanders and the coefficient becomes only weakly significant.

The estimations suggest that increasing of one unit the absolute difference  $Abs\_p0Belifs$  decreases for a bystander the probability of modifying punishment decision across periods of approximately 3%.

Therefore, we conclude that also this second set of results supports the Differential Social Influence hypothesis

- Conclusion relative to the Differential Social Influence hypothesis: *There is evidence that the more a bystander punishes in the first period, the more she is responsive to the social information received. We also find evidence that the larger is the absolute difference between a bystander's beliefs regarding peers' average punishment and her first period punishment, the less likely is that she modifies punishment decisions across periods.*

*Therefore, we conclude that the results of our experiment support the Differential Social Influence hypothesis.*

#### 4.3. Equivalence of Normative and Informational Influence hypothesis

We conclude this section reporting results on the difference between bystanders exposed to both *Normative* and *Informational* social influence and those exposed only to the latter. From the summary statistics reported in Table 1, we could see that in the first period third-parties in *Informational* punish on average 4.1 tokens versus 3.3 of those in *Normative*. This difference is not statistically significant (p-value 26%). In both treatments, on average bystanders reduce punishment between first and second period: *Normative* of 0.28 tokens and *Informational* of 0.33. Also this difference is not statistically significant.

We test the hypothesis that third-parties in *Normative* are more likely to revise their second period punishment decisions. We create the dummy variable *Normative* equal to 1 for third-parties in the *Normative* treatment and we implement a logistic regression. The dependent variable is *DummyP1p0* equal to 1 when punishment is modified across periods. Results are reported in Table 4.3.

Table 6: Probability Modify Punishment Across Periods: Treated Groups

	(1)	(2)	(3)	(4)
<i>Normative</i>	0.064 (0.11)	0.075 (0.09)	0.034 (0.17)	0.154 (0.15)
Strat_Punish			0.054*** (0.01)	0.039** (0.02)
Abs_p0Belifs			-0.085*** (0.02)	-0.091*** (0.03)
NorStratPunish			0.016 (0.02)	0.014 (0.02)
NorAbs_p0Belifs			0.103*** (0.03)	0.116*** (0.04)
Abs_Signalp0			-0.003 (0.02)	0.001 (0.02)
NorAbs_Signalp0			0.007 (0.02)	-0.003 (0.02)
Abs_BelAvgPun			0.094*** (0.02)	0.094*** (0.03)
NormAbs_BelAvgPun			-0.104*** (0.03)	-0.117*** (0.04)
DummySignalp0			0.064 (0.10)	0.152 (0.11)
NorDummySignalp0			0.095 (0.14)	-0.038 (0.17)
Other contr	Y	Y	Y	Y
<i>N</i>	609	420	609	420
pseudo <i>R</i> <sup>2</sup>	0.092	0.092	0.316	0.221

Notes: Logistic regression: dep. var. *DummyP1p0*, marginal effect at means, SE clustered by subject. Significance levels: \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. Other *Controls* include: male age degree worker social arts field.other risk logic impulsivity Instruction DictatorTake.

From the coefficient of *Normative* in model 1 to 4 we could see that, on average, there is no statistical difference between treatments in the likelihood of modifying punishment decision. Moreover, we test if participants in *Normative* and *Informational* choose a second period punishment that

reduces the gap between first period individual punishment and peers' punishment with the same probability. We implement a logistic regression and we report the results in Table 5. From the coefficient of the dummy variable *Normative* in models (2) and (3) we could see that there is no statistical difference between the two effect treatments.

However, when instead we disentangle the effect of individual determinants of the probability to modify punishment decision across periods it is possible to find differences between treatments. First, consider the tests we did for the Differential Social Influence hypothesis. Models 7 and 8 of Table 3 suggest that increasing one unit the absolute difference across punishment in the first period and individual beliefs regarding the average punishment in the session (the variable *Abs\_p0Beliefs*) decreases of approximately 3% the likelihood to modify punishment between period. However, the result is only weakly significant. Nevertheless, the estimation could be affected by the fact that in the models of table 3 we constrained the slope of *Abs\_p0Beliefs* to be the same for *Normative* and INFORMATIONAL. Therefore, in model 3 and 4 of Table 4.3 we introduce the interaction term *NorAbs\_p0Beliefs*, that isolate the effects of the absolute difference between punishment in the first period and beliefs about peers' average punishment for bystanders in the *Normative* treatment. The coefficient is positive and significant at 1% level for both model specification, and the coefficient of *Abs\_p0Beliefs* becomes negative and significant at 1% level. Interpreting the coefficients, we can see that for third-parties in *Normative* *Abs\_p0Beliefs* has no effect on the probability of modifying punishment across periods. Instead, for bystanders in *Informational* an increase of one unit in *Abs\_p0Beliefs* diminishes of roughly 8% the probability of modifying punishment across periods.

We can find a similar difference between treatments for the probability that subjects choose a second period punishment that decreases the gap between individual first period punishment and peers' average punishment. From model (3) in Table 5 we can see that the unconstrained coefficient of



$Abs_{p0}Belifs$  is negative and statistically significant at the 1% level. However, the coefficient of the variable  $NorAbs_{p0}Belifs$  is positive and statistically significant. These results confirm that  $Abs_{p0}Belifs$  has a statistically significant effect only for subjects in the *Informational* treatment. Therefore, the results contrast the Differential Social Influence hypothesis for bystanders in the *Normative* treatment, while the hypothesis finds support for subjects in the *Informational* treatment.

As a possible explanation for this difference across treatments, we conjecture that for bystanders in *Normative* there are additional incentives to revise their punishment decisions compared to bystanders in INFORMATIONAL. In fact, in *Normative* bystanders are told that their punishment choices of the second period will be revealed to other participants and that these peers will express their judgements. Therefore, it seems that the threat of revealing individual choices to other participants triggers the decision to modify first period punishment. In order to further investigate this conjecture, in models (4) and (5) of Table 5 we restrict the analysis to those punishment choices within the TREATED group where  $Abs_{p0}Belifs$  is larger than the average. In fact, participants choosing a first period punishment largely different from their beliefs about average punishment show that they are not interested in conforming to peers' behavior and so they should be unaffected by *Informational* social influence. On the other hand, subjects in *Normative* before stating the second punishment decision are told that peers will observe and judge their individual choices. Therefore, this additional exposure to *Normative* social influence might induce some of the subjects unresponsive to *Informational* social influence to conform with the peers' punishment behavior. Indeed, in models (4) and (5) the coefficient of the dummy *Normative* is positive and statistically significant. These estimations suggest that, within the subsample of subjects showing lack of interest for conformity with peers' punishment behavior, participants exposed to *Normative* social influence are roughly 20% more likely to choose a second period punishment

that reduces the gap between first period individual punishment and peers' average punishment.

Furthermore, consider the absolute difference between a bystander's beliefs regarding peers' average punishment and the actual peers' average punishment. Models 7 and 8 of Table 3 indicate that increasing of one unit the coefficient of the variable *Abs\_BeliefAvgPunish* increases for a bystander the probability to modify punishment decision across periods by 3%. However, this result comes from models where we constrained the coefficient of *Abs\_BeliefAvgPunish* to be the same in *Normative* and *INFORMATIONAL*.

We verify if the coefficient is the same in both treatments estimating the effect of *Abs\_BeliefAvgPunish* for the two groups separately. We do so interacting the variable *Abs\_BeliefAvgPunish* with the dummy *Normative* and so creating the variable *NormAbs\_BeliefAvgPunish*. From results of models 3 and 4 in Table 4.3, we can see that the coefficient of *NormAbs\_BeliefAvgPunish* is negative and statistical significant at 1% level. On the other hand, the coefficient of *Abs\_BeliefAvgPunish* in the unconstrained model becomes positive and significant at the 1% level, while it was only weakly statistically significant in the constrained model. Specifically, for subjects in the *Informational* treatment an increase of one unit in *Abs\_BeliefAvgPunish* raises by approximately 9% the probability of modifying punishment across periods.

In order to further investigate this result, we check how bystanders in the two treatments modify their punishment choices across periods conditional to the sign of the difference between beliefs regarding peers' average punishment and actual peers' average punishment. Table 4.1 reports these summary statistics. Bystanders in both treatments reduce punishment in the second period when actual average punishment is lower than expected. However, bystanders in *Informational* on average reduce punishment of more than 1 tokens, while those in *Normative* of less than 0.5. Conversely, when actual average peers' punishment exceeds a bystander's expectations, in

*Informational* third-parties increase punishment of 0.4 tokens on average, while bystanders *Normative* do not modify punishment decisions.

It is possible that the lower variability registered in *Normative* derives from the fact that individual choices are observable by peers. We conjecture that in *Normative* bystanders refrain from modifying punishment decisions, in particular from reducing punishment, because of the disutility of being eventually judged and targeted with the "sad" emoticon by peers.

Finally, we also test if the slope of the variable *Strat\_Punish* differs *Normative* and INFORMATIONAL. We created the variables *NorStrat\_Punish*, isolating the effect of *Strat\_Punish* for bystanders in the *Normative* treatment. As expected, results for these unconstrained models reported in Table 4.3 show that there is no statistical difference between treatments in the data.

- Conclusion relative to Equivalence of *Normative* and *Informational* Influence hypothesis: *We find mixed evidence regarding our hypothesis. On one hand, at an aggregate level the likelihood to modify punishment choices is the same in Normative and INFORMATIONAL. However, disentangling the determinants that push bystanders to modify punishment choices across periods, we find differences between the two treatments. Therefore, we conclude that the empirical evidence is mixed and the hypothesis is not fully supported by the data.*

## 5. Conclusions

Human organizations need mechanisms to enforce rules and regulations upon which are founded. On one hand, societies developed a centralized apparatus of enforcement for this purpose. However this centralized systems coexist with a decentralized practice of punishment carried on by members of the societies itself. Understanding the nature and characteristics of decentralized punishment might help legal scholars and policymakers to design effective

policies in a variety of situations. Therefore, which are the major drivers of decentralized third-party punishment is an important question for social scientists.

In this paper we examine through a laboratory experiment the effect of one of these drivers, social influence, on the punishment decisions of third parties. Scholars in psychology, law and economics underline the relevance of third party punishment for the cohesion of human societies (Fehr and Fischbacher, 2004b; Marlowe et al., 2008) and the importance of social influence in various fields of application (Bernheim, 1994; Turner, 1991; Kahan, 1997; Becker, 1991). However, to the best of our knowledge this paper is the first work that investigates empirically social influence effects within the framework of third party punishment.

In a modified dictator game, we elicit the punishment choices of third parties before and after having exposed them to information regarding the punishment behavior of their peers. We compare those choices with decisions made by bystanders not exposed to social relevant information. The main finding of this article is that social influence is a major driver of bystanders' decision to engage in third-party punishment. Results of the experiment show that third-parties receiving information about peers' punishment revise their punishment choices more often and on average reduce punishment across periods compared to bystanders exposed to social irrelevant information. This last effect seems to be driven by the fact that bystanders' beliefs regarding peers' average punishment are higher than the actual punishment peers exert. Indeed, consistently with the model predictions, the empirical analysis shows that the larger is the absolute difference between a bystander's beliefs about peers' average punishment and peers' actual punishment, the more likely is the bystander to revise his initial decisions.

We also disentangled the effect of two possible channels of social influence. Results suggest that some third-parties are only responsive to the discomfort for disagreeing with the majority, that is at the base of *Normative* social

influence and their punishment choices are not influenced by the "need to be right" on which is based *Informational* social influence. Distinguishing between these two channels of social influence is of primary importance for social analysts, since previous studies document that *Informational* social influence causes a permanent change in behavior (see for example Newcomb et al., 1967). On the other hand, *Normative* social influence is more ephemeral and leads to modifications of behavior that are subject to specific circumstances<sup>17</sup> (Deutsch and Gerard, 1955; Cohen and Golden, 1972; Burnkrant and Cousineau, 1975).

These findings have two major implications. On one hand, they remark the importance in our societies of citizens' perception about peers' behavior. This is especially important in situations where beliefs of the general population systematically overestimate the frequency of socially undesirable behaviors, like frequently happens for perceived crime, benefit frauds or percentage of non-voters<sup>18</sup>. In these situations, policymakers might often achieve welfare-improving results by means of ad-hoc communication strategies that could overperform alternative and often more costly policies (see for example Casal and Mittone, 2014, where the authors discuss an application of social stigma to tax evasion).

On the other hand, even when population beliefs are not biased, the possibility of resorting to social influence as a subsidiary tool for achieving compliance has been advanced by scholars in an array of situations of economic importance (Ela, 2008; Posner, 2000; Cooter, 1998; Zasu, 2007). As a society, we invest considerable amount of resources with the objective of shaping

---

<sup>17</sup>Nevertheless, scholars proposed models of endogenous preferences, arguing that even individuals initially adopting compliant behaviors by means of *Normative* social influence may endogenously modify their preferences (Akerlof, 1989; Klick and Parisi, 2008).

<sup>18</sup>For example, the Royal Statistical Society reports that 58% of the UK population estimates that crime is rising, while data show how crime rate in the country is 19% lower than the previous year and 53% lower than 1995. For discussion of other examples and additional details see <http://www.kcl.ac.uk/newsevents/news/newsrecords/2013/07-July/Perceptions-are-not-reality-the-top-10-we-get-wrong.aspx>.

individual beliefs and direct them toward social desirable outcomes. Policy-makers might want to encourage, by means of a social influence approach, third party interventions in situations where a centralized sanctioning authority might lack the ability or the resources to result effective. This is the case for example of the recent campaign aiming at prevention of social offenses "*Bringing in the Bystander*" promoted in the UK by the National Sexual Violence Resource Center. The campaign aims at reducing social offenses employing a marketing campaign that explicitly encourages third parties intervention<sup>19</sup>.

We agree that third-party punishment plays a key role in cementing human societies together (Mathew and Boyd, 2011). In this article we argue for the first time about the possibility for policymakers to take advantage of social influence effects in promoting third-party punishment, reporting evidence from a laboratory experiment that social influence significantly affects bystanders' interventions. Given the importance and wide possibilities of application in the societal framework, we hope that future researches further investigate the connection between social influence and third party punishment, in particular verifying the robustness of our findings in a field setting and the persistence of the effects in a longer term horizon.

---

<sup>19</sup>"Using a bystander intervention approach combined with a research component, this program assumes that everyone has a role to play in prevention [...] The Know Your Power campaign is the social marketing component of *Bringing in the Bystander*".

Appendix A.

First Period Punishment

Treatment	StratP0	StratP5	StratP10	StratP15	StratP20	StratP25	StratP30
<i>Control</i>							
(mean)	1.18	2.07	2.40	2.96	4.13	4.82	5.79
(median)	0	0	2	3	4	5	5
(SD)	3.31	4.09	2.74	2.97	3.94	4.68	5.72
<i>Normative</i>							
	.36	1.73	2.58	3.58	4.29	4.87	5.73
	0	1	2	2	4	5	5
	1.05	2.86	3.13	4.36	4.44	5.03	6.01
<i>Informational</i>							
	1.38	1.90	3.36	3.74	5.14	6.17	7.33
	0	1	3	4	5	6.5	8.5
	3.60	2.16	3.83	3.63	5.00	5.48	6.15
Total	.96	1.90	2.77	3.42	4.51	5.27	6.26
	0	1	2	3	5	5	6
	2.88	3.14	3.25	3.69	4.46	5.06	5.96

Table A.7: Average First Period Punishment by Levels of dictator Taking. StratP0 = dictator take 0 tokens from receiver

Second Period Punishment

Treatment	Punish0	Punish5	Punish10	Punish15	Punish20	Punish25	Punish30
<i>Control</i>							
(mean)	1.02	2.02	2.73	3.36	4.42	5.02	6.22
(median)	0	1	2	2	4	4	6
(SD)	3.18	3.45	3.49	3.58	4.38	5.24	6.26
<i>Normative</i>							
	.62	1.69	2.47	3.18	3.87	4.2	5.18
	0	1	2	3	3	4	4
	1.99	2.86	3.27	3.63	4.31	4.83	6.02
<i>Informational</i>							
	1.5	2.26	2.95	3.60	4.60	5.48	6.31
	0	1	3	4	5	6	8
	4.07	3.92	3.66	3.87	4.24	4.83	5.54
<b>Total</b>							
	1.04	1.98	2.71	3.37	4.29	4.89	5.89
	0	1	2	3	4.5	5	6
	3.16	3.41	3.45	3.67	4.29	4.96	5.93

Table A.8: Average Second Period Punishment by Levels of dictator Taking: Punish0 = dictator take 0 tokens from receiver



## Appendix B. Instructions

Welcome! This is a study on individual decision-making. Participants' answers are completely anonymous. It will not be possible to data analysts to link individual answers to the participants that provided them. You earned euro 5 for showing up on time today. Additionally, you can collect other earnings. The amount of these earnings depends on your choices and from the choices other participants will make during the study. During the study you will earn "tokens". Each 10 tokens earned, euro 1 will be paid out to you. In the unlikely case you will collect negative earnings, those losses will be subtracted from your participation fee. If you have questions at any time, please raise your hand and wait for a researcher that will answer your questions privately. Please switch off and remove from the table any electronic device, do not talk or communicate with other participants during the study. The study is composed of more parts. Earnings obtained in each part of the study are independent from those obtained in the other parts. Your final earnings are composed by:

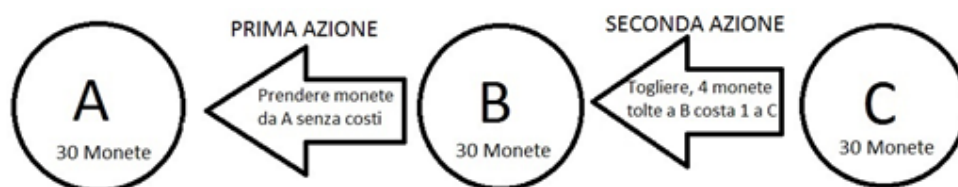
- Euro 5 of the participation fee
- Earnings collected in the first part of the study
- Earnings collected in one part after the first one. At the end of the study the computer will randomly select the part after the first one your earnings will be paid out to you

Final earnings will be paid privately and cash at the end of the study

First Part Instructions: description of the situation (Instructions on this part are the same in the 3 treatments)

Consider a situation with 3 people. Each person is randomly assigned to a role: one "Person A", one "Person B" and one "Person C". A, B and C could make decisions and earn tokens.

- Person A receives 30 tokens and does not make decisions
- Person B receives 30 tokens. Moreover, B could take some or all A's tokens and add them to his own earnings without incurring costs. Precisely, B could take 0, 5, 10, 15, 20, 25 or 30 tokens from A.
- Person C receives 30 tokens, observes B's action and could eliminate some of B's tokens, incurring a cost. For each 4 tokens eliminated from B's earnings, A has to pay 1 token. Person C could use up to 20 tokens to reduce B's earnings. C's decision does not affect A's earnings



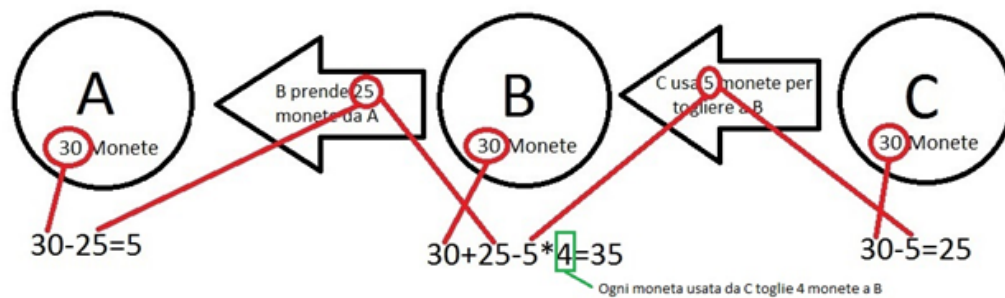
Therefore, A, B and C earnings are:

- Person A:  $(30 \text{ initial tokens}) - (\text{tokens taken by B})$
- Person B:  $(30 \text{ initial tokens}) + (\text{tokens taken from A}) - (4 * \text{tokens used by C})$
- Person C:  $(30 \text{ initial tokens}) - (\text{tokens used for reducing B's earnings})$

Example 1) (please look at your computer screen): B takes 25 tokens from A. After observing B's choice, C decides to use 5 tokens to reduce B's earnings. Therefore participants' final earnings are:

- Person A = 5 tokens (tokens left by B)

- Person B = 35 tokens (30 initial tokens + 25 tokens taken from A –  $5 \cdot 4 = 20$  tokens coming from the 5 tokens used by C to reduce B's earnings)
- Person C = 25 tokens (30 initial tokens – 5 tokens used to reduce B earnings)

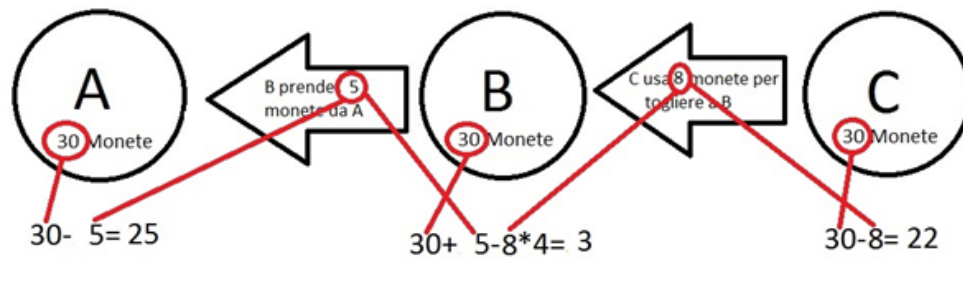


Example 2) (please look the computer screen): B takes 5 tokens from A. After observing B's choice, C uses 8 tokens to reduce B's earnings. Therefore participants' final earnings are:

- Person A = 25 tokens (left by B)
- Person B = 3 tokens (30 initial tokens + 5 tokens taken from A –  $8 \cdot 4 = 32$  tokens coming from the 8 tokens used by C to reduce B's earnings)
- Person C = 22 tokens (30 initial tokens – 8 tokens used to reduce B's earnings)

Your actions and earnings

Person C observes how many tokens B takes from A. You and the other participants in the laboratory have to indicate the number of tokens, an



integer between 0 and 20, that C in your opinion will use. When everyone have answered, we calculate the average of the individual amounts indicated by you and the other participants. If the number you indicated is equal, bigger or smaller by one unit respect the average, you receive 40 tokens that will be added to your final earnings (if you indicate 0, you will receive the forty tokens if the average is 0, 1 or 2; if you indicate 20, you will receive the 40 tokens if the average is 20, 19 or 18). Instead, you do not earn tokens in this part of the study if the number you indicate is bigger or smaller by more than one unit with respect to the average.

Example 1) (please look at the computer screen): Consider the action of B “take 20 tokens from A and collect 50 tokens, leaving 10 tokens to A”. You indicate that C uses 11 tokens. You receive 40 tokens if on average all the participants to the study indicated “11”, “10” or “12” tokens. If the average is different from these numbers, you will not earn tokens for this part of the study

Example 2) (please look at the computer screen): Consider the action of B “take 0 tokens from A and collect 30 tokens, leaving 30 tokens to A”. You indicate that C uses 3 tokens. You receive 40 tokens if on average all the participants to the study indicated “3”, “2” or “4” tokens. If the average is different from these numbers, you will not earn tokens for this part of the study.

Tempo rimanente per prendere una decisione 113

La tabella sotto elenca le possibili scelte della Persona B. Quante monete uscirà la Persona C? Guadagni 40 monete se la tua risposta è uguale oppure maggiore/minore di una moneta a quella media fornita dai partecipanti.

Scelta Persona B	Quante monete uscirà la Persona C?
Prendere monete 30 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 0, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 25 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 5, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 20 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 10, la Persona C monete 30 monete uscite)	<input type="text" value="11"/>
Prendere monete 15 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 15, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 10 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 20, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 5 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 25, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 0 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 30, la Persona C monete 30 monete uscite)	<input type="text"/>

[Continua](#)

**Guadagni le 40 monete se in media i partecipanti allo studio avranno indicato "11" oppure "10" oppure "12".  
Guadagni 0 se la media è diversa da questi valori.**

Tempo rimanente per prendere una decisione 293

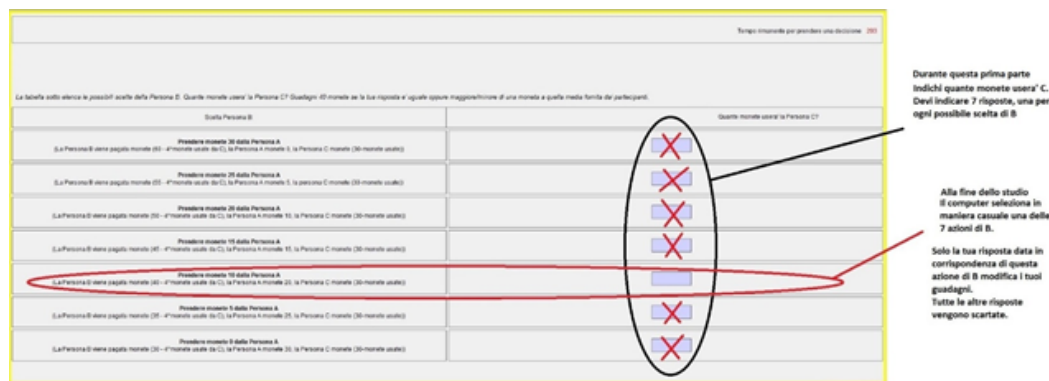
La tabella sotto elenca le possibili scelte della Persona B. Quante monete uscirà la Persona C? Guadagni 40 monete se la tua risposta è uguale oppure maggiore/minore di una moneta a quella media fornita dai partecipanti.

Scelta Persona B	Quante monete uscirà la Persona C?
Prendere monete 30 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 0, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 25 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 5, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 20 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 10, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 15 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 15, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 10 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 20, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 5 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 25, la Persona C monete 30 monete uscite)	<input type="text"/>
Prendere monete 0 della Persona A (La Persona B viene pagata monete 05 - 4 monete uscite da C, la Persona A monete 30, la Persona C monete 30 monete uscite)	<input type="text" value="3"/>

[Continua](#)

**Guadagni le 40 monete se in media i partecipanti allo studio avranno indicato "3", "2" oppure "4". Guadagni 0 se la media è diversa da questi valori**

You are required to indicate how many tokens Person C uses for each possible action of B (B takes 30 tokens from A; B takes 25 tokens. . . ; B takes 0 tokens from A). At the end of the study, the computer will randomly select one of the 7 actions of Person B. Relatively to this action, we will verify if you earned the 40 tokens. Your decisions and those of the other participants relative to other possible actions of Person B will be discarded and will not affect your final earnings.



Before starting this first part of the study, we ask you to answer some *Control* questions. Answers to these *Control* questions will not affect your final earnings.

(Participants answer *Control* questions on their computers. The Ztree file containing the *Control* questions is available upon request to the authors).

Instruction second part: description of the situation (instructions on this part are the same in all treatments)

Consider the same situation described in the first part, where 3 people are present, A, B and C, that can make decisions and earn tokens. Exactly as in the first part:

- A receives 30 tokens and do not makes decisions

- B receives 30 tokens and could take some or all the tokens of A
- C receives 30 tokens, observes the action of B and could reduce earnings of B paying a cost (every 4 tokens of reduction of B's earnings C has to pay 1 token)

Your actions and earnings

In this second part you and the other participants have to make decisions first as “Person B” then as a “Person C”. Therefore, you have to indicate:

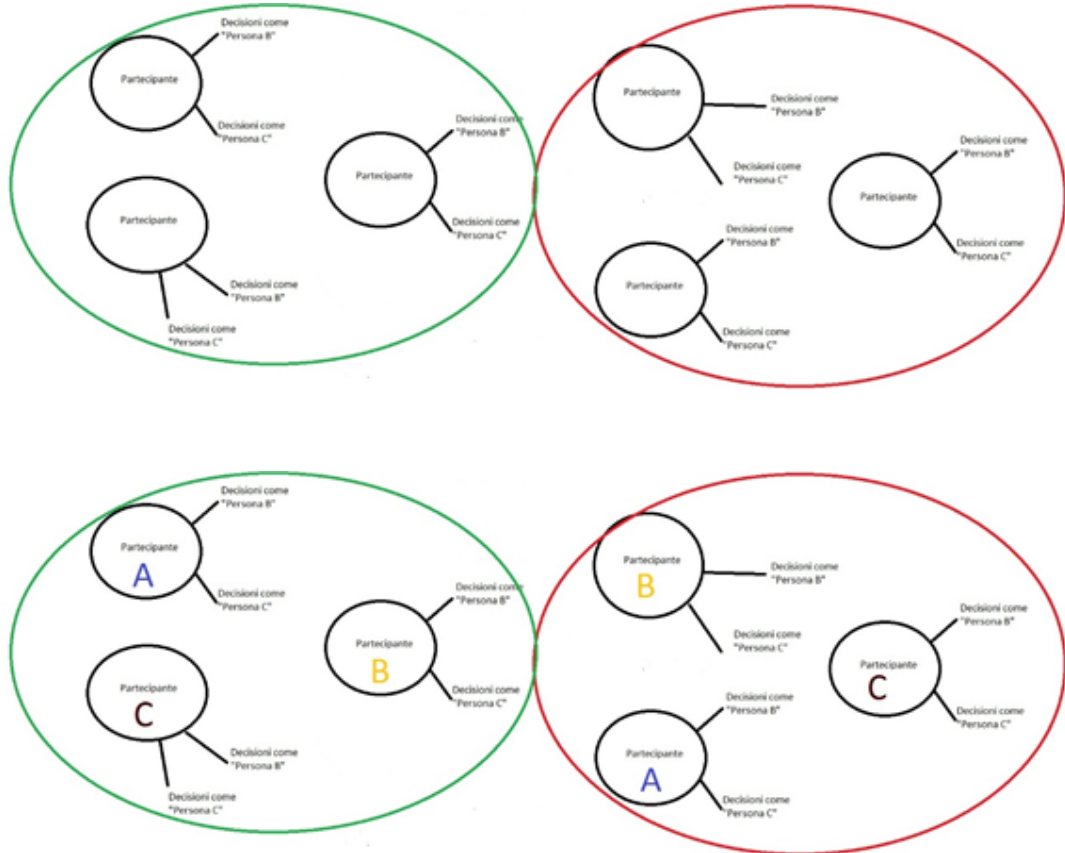
- First, as “Person B”, how many tokens you take from A
- After, as a “Person C”, for any possible action of B how many tokens you use for reducing B's earnings



Why do you have to make decisions both as “Person C” and as a “Person B”? In calculating final earnings, each participant is associated to an unique role: either Person A or Person B or Person C. However, you and the other participants will not know which role has been assigned to you until the end of the study today. Indeed, you and the other participants will be randomly divided in groups of 3

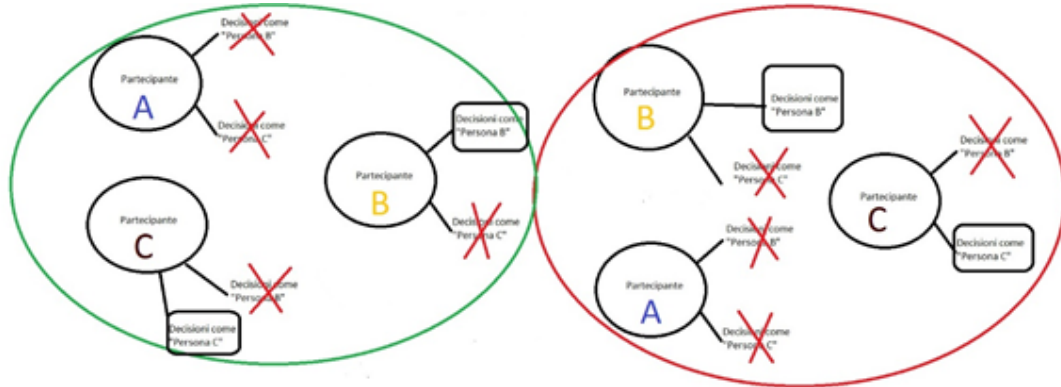
Within the group, each one of the 3 participants is assigned either to role A, B or C

Assignment to groups and assignment of roles is completely random and each participant has 1 possibility over 3 of being assigned a specific role.



Therefore, if you are assigned the role “Person A”, your final earnings are determined by the tokens left you by the Person B that is in your same group. Other decisions you make as a Person B or C will be discarded and have no influence on your final earnings nor on the earnings of the other participants. Similarly, participants assigned to the role “Person B” determine their final earnings and those of the other group components only by the decisions make as Person B. Decisions made as Person C have no effects on final earnings. Finally, also Participants assigned to role “Person C” only influence final earnings only by decisions make as C.





During this second part of the study we will also ask to indicate the day of the month in which you were born (E.g. if you were born January 25th 1983 you should report “25”).

Earnings of A, B and C in this second part are determined exactly as in the first part:

- Person A: (30 initial tokens) – (tokens taken by B)
- Person B: (30 initial tokens) + (tokens taken from A) – (4\*tokens used by C)
- Person C: (30 initial tokens) – (tokens used for reducing B’s earnings)

Before starting this first part of the study, we ask you to answer some *Control* questions. Answers to these *Control* questions will not affect your final earnings.



Instruction third part (*Normative Treatment*; instructions for *Control* and *Informational* are available upon request)

Now it starts the third and last part of this study. After the end of this part, we will ask you to fill in a brief questionnaire and then we will proceed with payments. Consider exactly the same situation of the second part of the study, same roles of A, B and C, same possible decisions that B and C have to make and same initial endowments and possible earnings. As in the second part, you have to make decisions first as a Person B then as a Person C. Additionally, in this third part before making your decisions you will receive information regarding the other participants. You will receive information on decisions made as Person C by the participants at today study. You will know how many tokens on average participants used in the second part of the study to reduce B's earnings. You will receive this information for any of the 7 possible B's choices.

Furthermore, before the end of the study, individual decisions as “Person C” that you are going to make in this third part will be revealed to 5 participants randomly selected. Similarly, you will received information regarding the individual choices made as Person C by 5 other participants. Each participant will be randomly assigned to an ID number. The ID number assigned is independent from the number of PC you are sit on. After you saw the individual choices of the other 5 participants, you and the other

Monete che B prende da A	Media monete riduzione guadagni B parte 2	Monete usate Partecipante 1 parte 3	Monete usate Partecipante 2 parte 3	Monete usate Partecipante 3 parte 3	Monete usate Partecipante 4 parte 3	Monete usate Partecipante 5 parte 3
0	0	0	0	0	0	0
5	0	0	0	0	0	0
10	0	0	0	0	0	0
15	0	0	0	0	0	0
20	0	0	0	0	0	0
25	0	0	0	0	0	0
30	0	0	0	0	0	0

participants will be able to vote for sending a smiling or sad emoticon

Partecipante 15	<input type="checkbox"/> Sorridente	<input type="checkbox"/> Triste
Partecipante 16	<input type="checkbox"/> Sorridente	<input type="checkbox"/> Triste
Partecipante 17	<input type="checkbox"/> Sorridente	<input type="checkbox"/> Triste

You receive a smiling emoticon if the majority of the five participants that saw your choices vote for “smiling”. Otherwise you will receive a sad emoticon. The emoticon will remain on your screen for one minute, then disappears automatically. After this minute has passed, you will know your final earnings.

If you have questions, please raise your hand and we will answer to you privately. Otherwise push “Continue” button and start with the third part.

## Appendix C. Description of the Variables Used in the Regressions


Table C.9: Variables


Variable	Description
degree	1 if subject completed 8th grade ("scuola media"), 2 if subject completed high school, 3 if subject has a bachelor degree or equivalent, 4 if subject has a master degree or equivalent, 5 if subject has a PhD or equivalent
worker	binomial variable, 1 if worker
male	binomial variable, 1 if male
age	subject's age
social	binomial variable, 1 if subject is a student in social sciences and medicine
arts	binomial variable, 1 if subject is a student in arts or humanities
field_other	binomial variable, 1 if subject not in social or arts
DictatorTake	total amount of tokes a subject when choosing as a dictator takes to the receiver in the 2 periods
risk	$\in [1, 10]$ , 1 if to question "In general, do you consider yourself ready to take risks?" the answer is "Not at all", 10 if the answer is "Totally ready to take risks"
logic	$\in [0, 2]$ , 1 point for each correct answer. See figures C.3 and C.4 below for the 2 questions.

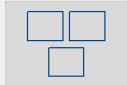
Variable	Description
impulsivity	$\in [0, 3]$ , 1 point for each correct answer. See figures C.5, C.6 and C.7 below for the 3 questions.
<i>Normative</i>	binomial variable, 1 for subjects in <i>Normative</i> treatment
TREATED	binomial variable, 1 for subject either in <i>Normative</i> or in <i>Informational</i> treatments
Strat_Punish	punishment first period
Beliefs_Punish	beliefs about peers' average punishment first period
Abs_p0Belifs	absolute value (Strat_Punish - Beliefs_Punish)
Abs Signalp0	absolute vale (Strat_Punish - Peers' average punishment period 1)
Abs BelAvgPun	absolute value (Beliefs_Punish - Peers' average punishment period 1)
NorStratPunish	<i>Normative</i> *Strat_Punish
NorAbs p0Belifs	<i>Normative</i> *Abs_p0Belifs
NormAbs_BeliefAvgPunish	<i>Normative</i> *Abs_BeliefAvgPunish
Diff_Deviation	abs (Strat_Punish - Peers' average punishment period 1) - abs (Punish - Peers' average punishment period 1)
ConvergeDummy	dummy, = 1 if Diff_Deviation>0
Instructions	total time employed by subjects to correctly answer <i>Control</i> questions

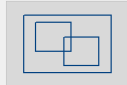
**DOMANDA 6**

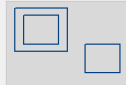
Quale tra questi diagrammi rappresenta la relazione tra:  
ARANCE-AGRUMI-FRUTTA  
(choose the answer and press OK)











**OK**

Figure C.3:

**DOMANDA 7**

Indica l'elemento grafico che completa la serie.  
(seleziona la risposta e premi OK)











**OK**

Figure C.4:

**DOMANDA 8**

Una mazza e una pallina costano euro 1.10 in total. La mazza costa euro 1.00 piu' della pallina.  
Quanto costa la pallina?

Euro

Figure C.5:

**DOMANDA 9**

Se 5 macchine impiegano 5 minuti per fare 5 oggetti,  
quanto tempo impiegheranno 100 macchine per fare 100 oggetti?

Minuti

Figure C.6:



In un lago c'è una chiazza di orchidee. Ogni giorno, questa chiazza raddoppia in dimensioni.  
Se la chiazza di orchidee impiega 48 giorni per coprire completamente il lago, quanto tempo impiegherà per coprire la metà del lago?

Giorni

Figure C.7:

## References

- Akerlof, G. A., 1980. A theory of social custom, of which unemployment may be one consequence. *The quarterly journal of economics* 94 (4), 749–775.
- Akerlof, G. A., 1989. The economic of illusion. *Economics & Politics* 1 (1), 1–15.
- Akerlof, G. A., Yellen, J. L., Katz, M. L., 1996. An analysis of out-of-wedlock childbearing in the united states. *The Quarterly Journal of Economics* 111 (2), 277–317.
- Asch, S., 1951. Effects of Group Pressure upon the Modification and Distortion of Judgments.
- Asch, S., 1956. Studies of Independence and Conformity: A minority of One Against a Unanimous Majority. *Psychological Monographs: General and Applied* 70 (9), 1–70.
- Balafoutas, L., Nikiforakis, N., 2012. Norm enforcement in the city: A natural field experiment. *European Economic Review* 56 (8), 1773–1785.
- Becker, G. S., 1991. A note on restaurant pricing and other examples of social influences on price. *Journal of Political Economy* 99 (5), 1109.
- Bernhard, H., Fischbacher, U., Fehr, E., 2006. Parochial Altruism in Humans. *Nature* 442 (7105), 912–915.
- Bernheim, B., 1994. A Theory of Conformity. *Journal of political Economy*, 841–877.
- Bond, R., Smith, P., 1996. Culture and Conformity: A Meta-analysis of Studies using Asch’s (1952b, 1956) Line Judgment Task. *Psychological Bulletin; Psychological Bulletin* 119 (1), 111.

- Buckholtz, J., Asplund, C., Dux, P., Zald, D., Gore, J., Jones, O., Marois, R., 2008. The Neural Correlates of Third-party Punishment. *Neuron* 60 (5), 930–940.
- Burnkrant, R. E., Cousineau, A., 1975. Informational and normative social influence in buyer behavior. *Journal of Consumer research*, 206–215.
- Casal, S., Mittone, L., 2014. Social esteem versus social stigma: the role of anonymity in an income reporting game. CEEL Working Papers 1401, Cognitive and Experimental Economics Laboratory, Department of Economics, University of Trento, Italia.
- Cason, T. N., Mui, V.-L., 1998. Social influence in the sequential dictator game. *Journal of Mathematical Psychology* 42 (2), 248–265.
- Chavez, A. K., Bicchieri, C., 2013. Third-party sanctioning and compensation behavior: Findings from the ultimatum game. *Journal of Economic Psychology* 39 (December), 268–277.
- Christakis, N. A., Fowler, J. H., 2007. The spread of obesity in a large social network over 32 years. *New England journal of medicine* 357 (4), 370–379.
- Cialdini, R. B., 1993. *Influence: the Psychology of Persuasion*. Collins, NY.
- Cialdini, R. B., Goldstein, N. J., 2004. Social influence: Compliance and conformity. *Annu. Rev. Psychol.* 55, 591–621.
- Coffman, L., 2011. Intermediation Reduces Punishment (and Reward). *American Economic Journal: Microeconomics* 3 (4), 77–106.
- Cohen, J. B., Golden, E., 1972. Informational social influence and product evaluation. *Journal of Applied Psychology* 56 (1), 54.
- Coleman, S., 1996. The minnesota income tax compliance experiment: State tax results.

- Cooper, D., Rege, M., 2008. Social Interaction Effects and Choice under Uncertainty: an Experimental Study. Tech. rep., University of Stavanger.
- Cooter, R., 1998. Expressive law and economics. *The Journal of Legal Studies* 27 (2), 585–608.
- Cox, J. C., Friedman, D., Gjerstad, S., 2007. A tractable model of reciprocity and fairness. *Games and Economic Behavior* 59 (1), 17–45.
- Deutsch, M., Gerard, H. B., 1955. A study of normative and informational social influences upon individual judgment. *The journal of abnormal and social psychology* 51 (3), 629.
- Devenow, A., Welch, I., 1996. Rational Herding in Financial Economics. *European Economic Review* 40 (3), 603–615.
- Ela, J. S., 2008. Law and norms in collective action: Maximizing social influence to minimize carbon emissions. *UCLA Journal of Environmental Law & Policy* 27 (1).
- Engel, C., 2011. Dictator games: a meta study. *Experimental Economics* 14 (4), 583–610.
- Falk, A., Fischbacher, U., 2002. “Crime” in the Lab-Detecting Social Interaction. *European Economic Review* 46 (4), 859–869.
- Falk, A., Fischbacher, U., Gächter, S., 2010. Living in Two Neighborhoods—Social Interaction Effects in the Laboratory. *Economic Inquiry*.
- Falk, A., Ichino, A., 2006. Clean Evidence on Peer Effects. *Journal of Labor Economics* 24 (1), 39–57.
- Fehr, E., Fischbacher, U., 2004a. Social norms and human cooperation. *Trends in cognitive sciences* 8 (4), 185–190.

- Fehr, E., Fischbacher, U., 2004b. Third-party Punishment and Social Norms. *Evolution and human behavior* 25 (2), 63–87.
- Fehr, E., Gächter, S., 2002. Altruistic Punishment in Humans. *Nature* 415, 137–140.
- Fischbacher, U., 2007. z-Tree: Zurich Toolbox for Ready-made Economic Experiments. *Experimental Economics* 10 (2), 171–178.
- Fortin, B., Lacroix, G., Villeval, M., 2007. Tax Evasion and Social Interactions. *Journal of Public Economics* 91 (11), 2089–2112.
- Galbiati, R., Zanella, G., 2012. The tax evasion social multiplier: Evidence from italy. *Journal of Public Economics* 96 (5), 485–494.
- Gintis, H., 2000. Beyond *homo economicus*: evidence from experimental economics. *Ecological economics* 35 (3), 311–322.
- Glaeser, E. L., Sacerdote, B., Scheinkman, J. A., 1996. Crime and social interactions. *The Quarterly Journal of Economics* 111 (2), 507–548.
- Greiner, B., 2004. An Online Recruitment System for Economic Experiments.
- Hirshleifer, D., Hong Teoh, S., 2003. Herd Behaviour and Cascading in Capital Markets: a Review and Synthesis. *European Financial Management* 9 (1), 25–66.
- Hoff, K., Kshetramade, M., Fehr, E., 2011. Caste and Punishment: the Legacy of Caste Culture in Norm Enforcement. *The Economic Journal* 121 (556), F449–F475.
- Kahan, D. M., 1997. Social influence, social meaning, and deterrence. *Virginia Law Review*, 349–395.

- Klick, J., Parisi, F., 2008. Social networks, self-denial, and median preferences: Conformity as an evolutionary strategy. *The Journal of Socio-Economics* 37 (4), 1319–1327.
- Krupka, E., Weber, R., 2009. The Focusing and Informational Effects of Norms on Pro-social Behavior. *Journal of Economic Psychology* 30 (3), 307–320.
- Krupka, E. L., Weber, R. A., 2013. Identifying social norms using coordination games: Why does dictator game sharing vary? *Journal of the European Economic Association* 11 (3), 495–524.
- Kurzban, R., DeScioli, P., O'Brien, E., 2007. Audience Effects on Moralistic Punishment. *Evolution and Human behavior* 28 (2), 75–84.
- Lewis, P., Ottone, S., Ponzano, F., 2011. Free-Riding on Altruistic Punishment? An Experimental Comparison of Third-Party. *Review of Law and Economics* 7 (1).
- Lieberman, D., Linke, L., 2007. The Effect of Social Category on Third-party Punishment. *Evolutionary Psychology* 5 (2), 289–305.
- List, J. A., 2007. On the interpretation of giving in dictator games. *Journal of Political Economy* 115 (3), 482–493.
- Lotz, S., Baumert, A., Schlösser, T., Gresser, F., Fetchenhauer, D., 2011. Individual Differences in Third-Party Interventions: How Justice Sensitivity Shapes Altruistic Punishment. *Negotiation and Conflict Management Research* 4 (4), 297–313.
- Manski, C., 2000. Economic Analysis of Social Interactions. *Journal of Economic Perspectives* 14 (3), 115–136.
- Marlowe, F., Berbesque, J., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J., Ensminger, J., Gurven, M., Gwako, E., Henrich, J., et al., 2008. More

- "Altruistic" Punishment in Larger Societies. *Proceedings of the Royal Society B: Biological Sciences* 275 (1634), 587–592.
- Mas, A., Moretti, E., 2009. Peers at work. *American Economic Review* 99 (1), 112–145.
- Mathew, S., Boyd, R., 2011. Punishment Sustains Large-scale Cooperation in Prestate Warfare. *Proceedings of the National Academy of Sciences* 108 (28), 11375–11380.
- Nelissen, R., 2008. The Price You Pay: Cost-dependent Reputation Effects of Altruistic Punishment. *Evolution and Human Behavior* 29 (4), 242–248.
- Nelissen, R., Zeelenberg, M., 2009. Moral Emotions as Determinants of Third-party Punishment: Anger, Guilt, and the Functions of Altruistic Sanctions. *Judgment and Decision Making* 4 (7), 543–553.
- Newcomb, T. M., Koenig, K. E., Flacks, R., Warwick, D. P., 1967. Persistence and change: Bennington College and its students after twenty-five years. Wiley New York.
- Perkins, H., Linkenbach, J. W., Lewis, M. A., Neighbors, C., 2010. Effectiveness of social norms media marketing in reducing drinking and driving: A statewide campaign. *Addictive behaviors* 35 (10), 866–874.
- Piazza, J., Bering, J., 2008. Concerns about Reputation via Gossip Promote Generous Allocations in an Economic Game. *Evolution and Human Behavior* 29 (3), 172–178.
- Posner, E. A., 2000. Law and social norms: The case of tax compliance. *Virginia Law Review*, 1781–1819.
- Scharfstein, D., Stein, J., 1990. Herd Behavior and Investment. *The American Economic Review*, 465–479.

- Shinada, M., Yamagishi, T., Ohmura, Y., 2004. False Friends are Worse than Bitter Enemies: “Altruistic” Punishment of In-group Members. *Evolution and Human Behavior* 25 (6), 379–393.
- Sunstein, C. R., Schkade, D., Ellman, L. M., Sawick, A., 2006. Are judges political?: an empirical analysis of the federal judiciary. Brookings Institution Press.
- Swope, K., Cadigan, J., Schmitt, P., Shupp, R., 2008. Social position and distributive justice: Experimental evidence. *Southern Economic Journal*, 811–818.
- Topa, G., 2001. Social Interactions, Local Spillovers and Unemployment. *The Review of Economic Studies* 68 (2), 261–295.
- Turner, J. C., 1991. *Social influence*. Thomson Brooks/Cole Publishing Co.
- Zasu, Y., 2007. Sanctions by Social Norms and the Law: Substitutes or Complements? *The Journal of Legal Studies* 36 (2), 379–396.