

# On Punishment and Well-being

Jordi Brandts\* and María Fernanda Rivas<sup>‡</sup>

June 2007

## Abstract

The existence of punishment opportunities has been shown to cause efficiency in public goods experiments to increase considerably. In this paper we ask whether punishment also has a downside in terms of process dissatisfaction. We conduct an experiment to study the conjecture that an environment with stronger punishment possibilities leads to higher material but lower subjective well-being. The more general motivation for our study stems from the notion that people's subjective well-being may be affected by the institutional environment they find themselves in. Our findings show that harsher punishment possibilities lead to significantly higher well-being, controlling for earnings and other relevant variables. People derive independent satisfaction from interacting under the protection of strong punishment possibilities. These results complement the evidence on the neural basis of altruistic punishment reported in de Quervain et al. (2004).

JEL Classification Numbers: C92, D60, H40

Keywords: Public Goods, Experiments, Well-being, Punishment

---

\*Institut d'Anàlisi Econòmica (CSIC), Barcelona.

<sup>†</sup>The authors thank the Spanish Ministry of Education and Science and the Barcelona Economics program of CREA for financial support.

<sup>‡</sup>International Doctorate in Economic Analysis, Universitat Autònoma de Barcelona, and Departamento de Economía, Uruguay.

# 1 Introduction

As is well-known many situations of social interaction can be envisioned as public-good-type games in which individuals have incentives to take a "free-ride" on others' contributions to the public good and spend their own resources on other individually higher-valued uses. A large stream of experimental papers has documented that people often contribute to the public good and sometimes with large amounts of money. Nevertheless, observed inefficiency levels are still large and increase with experience, often ending up at near-zero provision.

However, in many such environments with free riding incentives, people do not need to passively accept the free riding of others. There often exist punishment opportunities of some sort, the possibility of taking actions that impose costs on others. Experimentalists have studied whether this possibility has any effect on social interaction, particularly when punishing others is costly. In an early contribution, Ostrom et al. (1992) study behavior in a repeated common pool resource game with uncertain horizon under different conditions involving punishment, communication and non-binding agreements. They find that under some conditions punishment opportunities lead to higher contribution levels. However, the fact that in their design the duration of the interaction is uncertain makes it possible that people develop an individual reputation so that there are material incentives for cooperation and punishment.

Fehr and Gächter (2000) report results from a finitely repeated public good experiment with and without costly punishment opportunities in which cooperation and punishment can never be part of subgame-perfect equilibrium, if rationality and selfishness are common knowledge. They provide very convincing evidence that the existence of punishment opportunities leads to a large increase in contributions. In one interesting extension Masclet et al. (2003) used experimental methods to study the power of informal non-material sanctions in a public good game and found that monetary and non-monetary sanctions initially increase contributions by a similar amount. Over time, however, monetary sanctions lead to higher contributions than non-monetary ones.

We take this evidence on the effectiveness of punishment as our starting point. If one simply stopped questioning at this point, social environments with strong punishment possibilities would appear preferable to environments that lack such possibilities. However, we believe that to make a judgement about the desirability of effective punishment possibilities one more element needs to be taken into consideration: the effect of punishment on process satisfaction. The general motivation for our study is the notion that people's subjective well-being may be affected by

the institutional environment they find themselves in and that economists need to understand these relations. Rabin (1993) formulates this as follows: *"Welfare economics should be concerned not only with the efficient allocation of material goods, but also with designing institutions such that people are happy about the way they interact with others. (...) Armed with well-founded psychological assumptions, economists can start to address the nonmaterial benefits and costs of the free market and other institutions."*

What led us to conducting the experiments presented in this study is the suspicion that the presence of punishment possibilities might have a downside, i.e. the possibility of using repressive sanctions may lead to low subjective well-being due to an uneasiness about the environment in which participants are immersed. If this were true, then it would not be straightforward to make an overall welfare judgement on the goodness of the presence of punishment possibilities, since two counter-vailing forces would have to be somehow compared with each other. If, in contrast, the presence of punishment possibilities had no significant or even a positive effect, then the judgement on the different institutional arrangements would be more direct, since both factors would point in the same direction.

We use the notion subjective well-being similarly to Kahneman, Wakker and Sarin's (1997) notion of "experienced utility", which goes back to Bentham. These authors consider that subjective well-being (or experienced utility) is both measurable and empirically distinct from standard decision utility. A subjective view of utility recognizes that everybody has his own ideas about happiness and the good life and that observed behavior is an incomplete indicator for individual well-being. Applied to our environment, it may well be that people make use of punishment possibilities and even that this leads to higher material payoffs. However, from here one can not directly conclude that people experience higher subjective well-being in such an environment.

We measure subjective well-being through self-assessments of participants' satisfaction with the experience in the experiment. Our focus is on the comparison of subjective well-being across different treatments in which we vary the punishment possibilities. In making this comparison we control for possible determinants of subjective well-being other than process considerations.

The novel aspect of this paper is precisely the analysis of the effects of punishment on subjective well-being. We believe that understanding this relation is important for a better understanding of social interactions. Our work is related to the study by de Quervain et al. (2004) on the neural basis of altruistic punishment. In that study subjects' brains were scanned

while they learned about the defector's abuse of trust and determined the punishment. It was found that the fact of effectively punishing a defector produces a satisfactory impact on the brain. What we study is not the direct effect of punishing on the punisher, but the effect on individuals of one of the features of the environment in which they interact. In the next section we present some background material on the measurement of well-being and on issues related to process satisfaction. After that we present the experimental design and procedures and then the results.

## 2 Background and previous evidence

Kahneman, Diener and Schwartz (1999) provide a wealth of information about the importance of well-being. Recent overviews about research into happiness and well-being and its relation to economics is provided among others by Frey and Stutzer (2002), Krueger (2005), and McFadden (2005). Veenhoven (1993) —the author is the founder of the Journal of Happiness Studies— presents a study on happiness all over the world.

We briefly discuss some of the previous work to illustrate the kinds of issues that are being studied. Van Praag and Ferrer-i-Carbonell (2004) present the perhaps most exhaustive study of what the authors call satisfaction analysis, based on the responses to subjective questions of the following type: *How satisfied are you with your financial situation, job, health, life, etc. Please respond on a scale from "very bad" to "very good" or on a numerical scale from 1 to 7 or 1 to 10.*

The authors of this study assert that humans do often evaluate many aspects of their situations guided by the objective of changing their life. They argue that "the empirical practice and success of these questions constitute ample evidence that individuals are able and willing to express their satisfaction on a cardinal scale. If we assume those questions to be interpreted in approximately the same way by different respondents and we find that similar respondent give similar answers, this is ample evidence that (approximate) interpersonal comparison is possible." They discuss how to study financial, job, housing, health, leisure, and environment satisfaction —what they call domains satisfactions— as well as satisfaction with life as a whole as a weighted aggregate of the domain satisfactions. The methodology is then applied to job satisfaction for a British data set and to political satisfaction for a Dutch survey.<sup>1</sup>

---

<sup>1</sup>There are many other research papers studying happiness. To name but a few we have Blanchflower and

Frey and Stutzer (2002) present an extensive literature survey. They report that one important finding of the literature about happiness is the large influence of non-financial variables on self-reported satisfaction. Frey and Stutzer (2000, 2002) classify the determinants of happiness into three blocks. The first group refers to micro and macro economic factors, the second one relates to institutional conditions in an economy and society, and the third group of determinants includes personality and demographic factors. With respect to the economic determinants of happiness, Frey and Stutzer report that in most nations, the fact of belonging to upper *income* groups somehow implies higher subjective well-being than belonging to lower income groups. However, the relation seems to be non-linear, there is diminishing marginal utility with absolute income. There may be many different reasons for that, one of the most important is that individuals compare themselves to others. Another explanation is in terms of aspiration levels (Easterlin, 2001). In this view happiness is determined by the gap between aspiration and achievement, and increases in income and aspiration levels are closely connected. An important economic determinant of happiness is *unemployment*. Being unemployed is correlated with low levels of satisfaction, not having a job imposes a high non-pecuniary stress and unhappiness.

In relation with the second group of determinants —institutional conditions—, Frey and Stutzer argue that the more developed *direct democracy* is, the happier the citizens are. Finally, with respect to the third group of determinants —personality and demographic factors— they find that people over 60 are happier than people under 30, people with higher education report higher well-being, and couples with and without children are happier than singles, single parents and people living in collective households.

There are a few experimental papers studying issues of well-being. Charness and Grosskopf (2001) analyze the relation between the importance people attach to relative payoffs and happiness, motivated by the conjecture that those who are less happy may seek solace in obtaining higher material payoffs than others. The experiment consisted in subjects making choice in simple dictator-type (one-shot) decisions tasks and in filling out a happiness questionnaire. The results summary is that there is no strong general correlation between happiness and concern for relative payoffs, but that the willingness to lower another person's payoff below one's own (competitive preferences) seems correlated with unhappiness. Brandts et al. (2004) study the impact of competition on the well-being of experimental subjects. Their approach is somewhat different from that of Charness and Grosskopf (2001). The idea is not to measure people's

---

Oswald (2000); Di Tella, MacCulloch and Oswald (2001); Diener and Oishi (2000); Diener and Seligman (2004).

homegrown levels of happiness but to evaluate whether different experiences in the lab could lead to different levels of process satisfaction. They find that competition has an adverse effect on the disposition towards others of those on the long side of the market and leads to lower subjective well-being for subjects on the long side of the market in comparison with those on the short side and those not subject to competition, all this controlling for earnings and other relevant variables.

Note, and this will become clearer below, that it would have been difficult to carry out this kind of studies on the basis of field data alone, since in natural environments it would be very hard to find adequate data with the desired parallel variations in the punishment conditions. It would probably have been even harder to obtain the corresponding information about subjective well-being.

The experimental design is explained in the next section and the results are presented in section 3. Finally, in section 4 we conclude.

### 3 Experimental design and procedures

In our experiments subjects interacted in pairs in a 20-round public goods game with punishment possibilities; the finite horizon was common information.<sup>2</sup> After the 20 rounds they had to answer one simple question about process satisfaction. There were two treatment variables: soft vs. strong punishment and partners vs. strangers matching. The first distinction responds to what motivates our study, while the second distinction will allow us to compare our work to that of Fehr and Gächter (2000). The difference between the types of punishment is the "fine-to-fee" ratio, which describes by how much the punished subject's income is reduced relatively to the fee the punishing subject has to pay to inflict punishment. This gives rise to a 2x2 treatment design which is summarized in table 1.<sup>3</sup> These are the essential features of our design.

|                  |          | Type of punishment   |                    |
|------------------|----------|----------------------|--------------------|
|                  |          | Strong               | Soft               |
| Type of matching | Stranger | Strong P, Stranger M | Soft P, Stranger M |
|                  | Partner  | Strong P, Partner M  | Soft P, Partner M  |

Table 1: Treatments

<sup>2</sup>Using pairs is the simplest starting point for studying what we are interested in. Group size effects could be studied in subsequent work.

<sup>3</sup>M refers to Matching and P to Punishment, hereafter.

Subjects were coupled randomly by the computer, maintaining the anonymity of the interaction partner.<sup>4</sup> In the Stranger Matching the groups —pairs— were changed round to round while in the Partner Matching the groups remained the same for all the rounds. The experiment was conducted at the Universitat Autònoma de Barcelona with undergraduate students of a variety of faculties. Participants were recruited by public advertisements posted throughout the campus.

Each round had two parts. In the first part each participant was asked to divide 5 tokens between two accounts, a group account —called account A in the experiment— and a private account —account B. Tokens placed in account A yield an identical amount of money to both members of a pair. Tokens in account B yield money only to the subject in question. Table 1 gives the payoff schedule, in tokens, depending on the contributions to account A, where X is the contribution of the subject in question and Y is the contribution of his partner. The first value of the cells is his payoff and the second one his partner’s payoff. Apart from these payoffs participants received 3 euros for showing up. The payoffs were calculated with a marginal rate of transformation between the public and the private account of 0.75, to create sufficient tension between the Nash equilibrium and the Pareto efficient allocation.

| X\Y | 0           | 1           | 2           | 3           | 4           | 5           |
|-----|-------------|-------------|-------------|-------------|-------------|-------------|
| 0   | 5 , 5       | 5.75 , 4.75 | 6.5 , 4.5   | 7.25 , 4.25 | 8 , 4       | 8.75 , 3.75 |
| 1   | 4.75 , 5.75 | 5.5 , 5.5   | 6.25 , 5.25 | 7 , 5       | 7.75 , 4.75 | 8.5 , 4.5   |
| 2   | 4.5 , 6.5   | 5.25 , 6.25 | 6 , 6       | 6.75 , 5.75 | 7.5 , 5.5   | 8.25 , 5.25 |
| 3   | 4.25 , 7.25 | 5 , 7       | 5.75 , 6.75 | 6.5 , 6.5   | 7.25 , 6.25 | 8 , 6       |
| 4   | 4 , 8       | 4.75 , 7.75 | 5.5 , 7.5   | 6.25 , 7.25 | 7 , 7       | 7.75 , 6.75 |
| 5   | 3.75 , 8.75 | 4.5 , 8.5   | 5.25 , 8.25 | 6 , 8       | 6.75 , 7.75 | 7.5 , 7.5   |

Table 2: Payoffs

The first part of each round was the same across the four treatments. After the first part of each round participants saw on their screens their partner’s decision and both payoffs —calculated following table 2. The second part of each round differed across the strong vs. soft punishment variation. In this part of the rounds, participants had the opportunity to punish their partners at a certain cost. The complete punishment schedule for both types of punishment can be seen in table 3. The cost of the punishment was the same, what changed was the punishment applied, i.e. the amount that participants could subtract from their partners’ payoff. In the Strong Punishment Treatment (STP, hereafter) the fine-to-fee ratio was 4 and in

---

<sup>4</sup>The experiment was programmed and conducted with the software z-Tree (Fischbacher, 2007).

the Soft Punishment Treatment (SOP, hereafter) it was 1.6.<sup>5</sup>

| <b>Strong Punishment Treatment</b> |          |             |             |             |             |             |             |             |             |             |             |             |             |
|------------------------------------|----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <b>Tokens deducted</b>             | <b>0</b> | <b>0.50</b> | <b>1.00</b> | <b>1.50</b> | <b>2.00</b> | <b>2.50</b> | <b>3.00</b> | <b>3.50</b> | <b>4.00</b> | <b>4.50</b> | <b>5.00</b> | <b>5.50</b> | <b>6.00</b> |
| Cost to the punishing subject      | 0.00     | 0.125       | 0.250       | 0.375       | 0.500       | 0.625       | 0.750       | 0.875       | 1.000       | 1.125       | 1.250       | 1.375       | 1.500       |
| Punishment level                   | 0        | 1           | 2           | 3           | 4           | 5           | 6           | 7           | 8           | 9           | 10          | 11          | 12          |

| <b>Soft Punishment Treatment</b> |          |             |             |             |             |             |             |             |             |             |             |             |             |
|----------------------------------|----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <b>Tokens deducted</b>           | <b>0</b> | <b>0.20</b> | <b>0.40</b> | <b>0.60</b> | <b>0.80</b> | <b>1.00</b> | <b>1.20</b> | <b>1.40</b> | <b>1.60</b> | <b>1.80</b> | <b>2.00</b> | <b>2.20</b> | <b>2.40</b> |
| Cost to the punishing subject    | 0.00     | 0.125       | 0.250       | 0.375       | 0.500       | 0.625       | 0.750       | 0.875       | 1.000       | 1.125       | 1.250       | 1.375       | 1.500       |
| Punishment level                 | 0        | 1           | 2           | 3           | 4           | 5           | 6           | 7           | 8           | 9           | 10          | 11          | 12          |

Table 3: Punishment costs

After the subjects had decided on the punishment level, they saw on their screens the first part payoffs, the punishment their partner decided to inflict on them, the cost of deducing tokens from their partner’s payoffs and their final payoffs —for the current period— calculated as:

$Final\ earning = Initial\ payoff\ (from\ table\ 2) - Tokens\ deducted\ by\ the\ partner - Cost\ of\ deducing\ tokens\ from\ the\ partner's\ payoff\ (from\ table\ 3)$ . Then the first part of the new round began and things proceeded in the same way until the end of round 20.

After the 20 rounds participants saw on their screens a summary of the experiment: contributions, punishment, and payoffs for each period. In the second part of the experiment we obtained our measurement of well-being. Each participant had to separately respond to the question *How satisfied would you say that you are with the experiment?* The possible answers were : 1) Completely satisfied, 2) Very satisfied, 3) Rather satisfied, 4) Neither satisfied nor dissatisfied, 5) Rather dissatisfied, 6) Very dissatisfied, 7) Completely dissatisfied.

After all this, we distributed a questionnaire; subjects did not know this beforehand.<sup>6</sup> In the questionnaire they were asked to indicate to what extent they agreed —on a six degree scale— with some statements, some of them referring to the punishment and others referring to their answer to the question about the satisfaction with the experiment. After they had filled out the questionnaire subjects were privately paid. Their final payoff was the total number of tokens earned (the sum for all the periods) converted into euros (1 token = 0.10 euros) plus the show-up fee (3 euros).

---

<sup>5</sup>An alternative to the strong vs. soft punishment distinction could have been a distinction between a punishment and a no-punishment treatment. We prefer our design choice, because it keeps the two treatments more parallel in procedural terms.

<sup>6</sup>Our well-being question is similar to one of the questions asked by Charness and Grosskopf (2001) in their more extensive happiness questionnaires.



## 4 Results

The main focus of our work is on the results on subjective well-being which will be presented in section 4.4. Before that we present some results on public goods contributions, punishment behavior and earnings which will allow us to relate our work to that of others and to better understand the results on well-being. Much of this presentation will be kept at a descriptive level.

|                  |          | Type of punishment |      | Total |
|------------------|----------|--------------------|------|-------|
|                  |          | Strong             | Soft |       |
| Type of matching | Stranger | 22                 | 22   | 44    |
|                  | Partner  | 24                 | 26   | 50    |
| Total            |          | 46                 | 48   | 94    |

Table 4: Number of subjects in the treatments

Tables 4 and 5 show some preliminary information. The most salient feature of the table 5 data is that the main difference is not between soft and strong punishment, but between partners and strangers. We did not foresee this. Observe also that contributions to the group account and the punishment applied are somewhat higher under strong punishment, for both strangers and partners, but that final earnings are lower under strong than under soft punishment. Finally, note that for partners, stronger punishment leads to somewhat higher final earnings; however, this is not true for strangers.

|                  |          |                               | Type of Punishment |        | Total |
|------------------|----------|-------------------------------|--------------------|--------|-------|
|                  |          |                               | Soft               | Strong |       |
| Type of Matching | Stranger | Contribution to group account | 1.14               | 1.32   | 1.23  |
|                  |          | Punishment applied            | 0.09               | 0.60   | 0.35  |
|                  |          | Final earnings                | 10.84              | 9.82   | 10.33 |
|                  | Partner  | Contribution to group account | 2.73               | 3.08   | 2.90  |
|                  |          | Punishment applied            | 0.30               | 0.33   | 0.31  |
|                  |          | Final earnings                | 11.76              | 12.25  | 12.00 |
| Total            |          | Contribution to group account | 2.00               | 2.24   | 2.12  |
|                  |          | Punishment applied            | 0.20               | 0.46   | 0.33  |
|                  |          | Final earnings                | 11.34              | 11.09  | 11.22 |

Table 5: Average values of the main results

We now study the behavior of the three variables shown in table 5 more in detail. Section 4.1 deals with contributions, 4.2 with punishment and 4.3 with final earnings. In each section we will summarize the main results as regularities.

## 4.1 Contributions to the group account

Figure 1 shows the average contributions per period —*per group of two*— for the four different treatments.

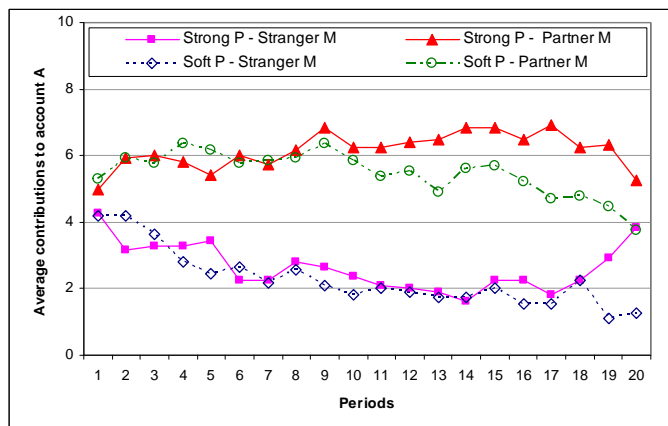


Figure 1: Contributions to the group account by type of P and M

The data presentation in the figure confirms the impression obtained from inspecting the averages in table 5. The main difference is due to the type of matching (partner or stranger) rather than the type of punishment (soft or strong). We do observe some secondary differences in that for both types of matching the contributions are higher for the Strong punishment case than the Soft punishment case, with larger differences for the Partner case, mostly in the last ten periods.<sup>7</sup> In the very last period, the average contribution in the SOP is 1.31, while in the STP it is 2.28. Observe also that in the Partner M case there is a tendency to diminish the contributions in the last periods while the tendency is the opposite in the STP-Stranger M.

Figure 2 gives a more aggregated view of the contribution data and makes the difference between the effects of the two treatments easier to see. It shows the contributions by type of punishment and by type of matching by group. The left panel of figure 2 shows that the average contribution to the group account is higher in the STP in comparison with the SOP, after round 8. While the average contribution in the first case remains almost the same, in the latter it decreases somewhat over time. If we look at the percentages of subjects that contributed zero tokens to the group account we find that they are 31.3 for the SOP and 13.0 for the STP. In

<sup>7</sup>With the exception of the last 2 periods, when the contribution of the STP-Stranger M increases substantially compared to the SOP-Stranger M, and in the final period it reaches the same value as the SOP-Partner M.

the figure on the right we see that the contribution in the Partner M Treatment is higher than in the Stranger M, although in the final periods they show a tendency to converge.

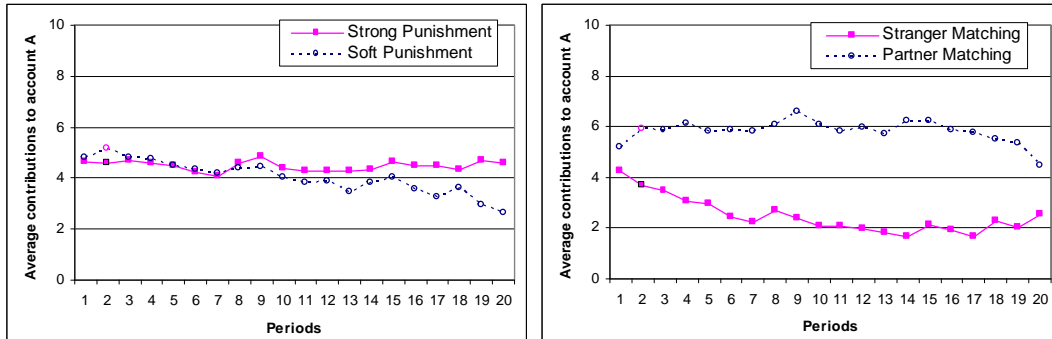


Figure 2: Contributions to the group account over time by type of P(left) and by type of M (right)

We can summarize these results in our first regularity.

*Regularity 1: Whether punishment is strong or soft has only secondary effects on contribution levels, the main difference is due to the type of matching.*

In their experiment Fehr and Gächter (2000) compared a situation with punishment with one without it. In both treatments, Stranger and Partner, they found that the existence of punishment rises contributions, and in the no punishment condition the contributions converge to full free-riding. Moreover, in their Partner treatment, the punishment opportunity makes contributions converge toward full cooperation.<sup>8</sup> Our results are different but not at odds with theirs. As can be seen in figure 2, the direction of the difference in contribution levels is in favor of strong punishment. It is just that the difference in contribution levels is small. The direction of the difference between strangers and partners is also the same direction as in Fehr and Gächter (2000). The difference in magnitudes can be explained by the difference in group size and marginal per capita return; we use  $n=2$  and  $MPCR=0.75$ , in contrast to their choices of  $n=4$  and  $MPCR=0.4$ .

<sup>8</sup>Fehr and Gächter (2000) also study the order effects of moving from a punishment to a no punishment environment and the reverse.

## 4.2 Punishment behavior

The fact that the punishment opportunities were different depending on the type of punishment makes it difficult to compare average punishment *amounts* between the two cases. What can be more easily compared is the percentage of subjects that punished their partner and this is what is done in table 6 and figures 3 and 4. In the table we see that the percentage of subjects that punished their partner at least once is higher in the STP, especially in the Stranger M case.

|                    | With soft punishment | With strong punishment |
|--------------------|----------------------|------------------------|
| Stranger Treatment | 55%                  | 91%                    |
| Partner Treatment  | 62%                  | 79%                    |
| Total              | 58%                  | 85%                    |

Table 6: Percentage of times subjects punished their partners at least once

In figure 3 we show the percentage of subjects that punished the partner by period for each one of the four treatments. In the Stranger M, the percentage is higher for the STP than for the SOP, in 18 out of the 20 periods (in the other 2 periods the percentages are the same). In the Partner M, in the first 12 periods the percentages are quite similar, but in the last 8 periods the percentage of subjects that punished the partner is higher in the SOP than in the STP.

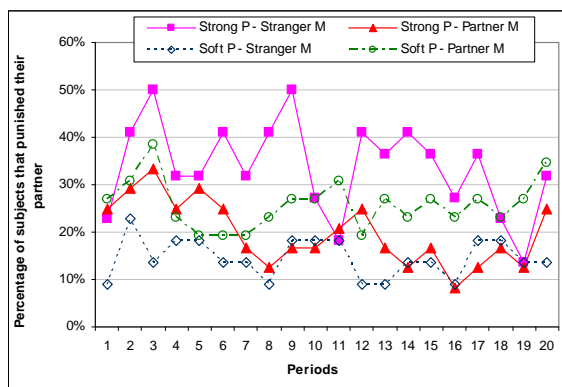


Figure 3: Percentage of subjects that punished their partner by type of P and M

In figure 4 on the left, we can see that the percentage is higher in the STP in the first periods, and in the last 5 they do not differ much —even in period 19 the percentage is higher for the SOP. In the figure on the right it can be seen that the evolution of the percentages are somewhat similar in both types of matching. Overall, the treatment differences in the use of punishment are not large.

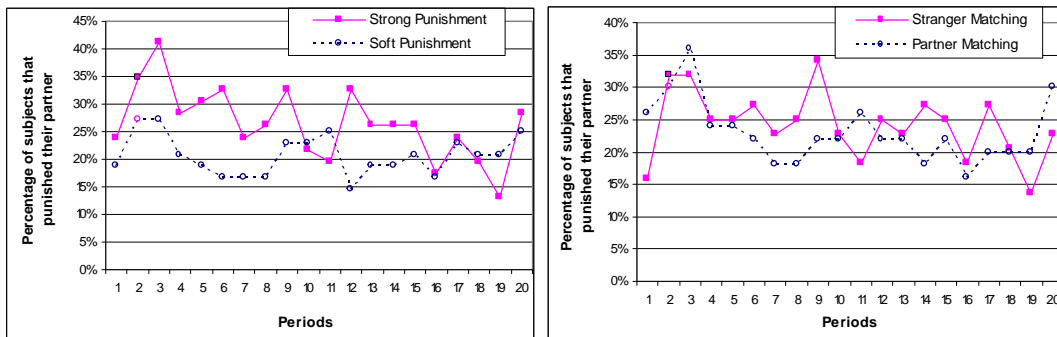


Figure 4: Percentage of subjects that punished their partner by type of P (left) and by type of M (right)

Figure 5 shows a histogram of punishment in terms of the punishment levels of table 3 for the four treatments. The key feature is that punishment is effectively not used very much, with the frequency of a zero level punishment being between 66% and 85%, depending on the treatment.

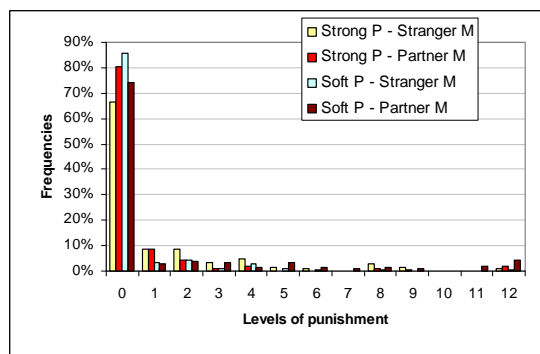


Figure 5: Histogram of punishment levels

To further understand the punishment process we want to relate the use of punishment not only to the treatment variables but also to behavior in earlier periods, in particular to the difference between a participant’s own contribution and the contribution of the partner. Given the high frequency of zeros in the levels of punishment, we decided to estimate a probit and a logit model where the dependent variable is a dummy that has value 1 if the subject punished his partner —regardless of with which amount— and has value 0 if the subject did not punish.<sup>9</sup>

<sup>9</sup>The difference between the ordered logit and the ordered probit model lies in the assumed distribution of  $\varepsilon_i$ . An ordered logit model assumes that  $\varepsilon_i$  is logistically distributed, while an ordered probit model assumes that it is normally distributed. The logistic distribution is similar to the normal except in the tails, that are heavier. See, for example, Wooldridge (2002).

To account for multiple observations in the estimation we clustered on subjects. In table 7 we present the estimates. The first three variables refer to the treatments. The variable *Neg\_deviat* (*Pos\_deviat*) represents the negative (positive) deviation of the partner's contribution from the subject's one. They are equal to:

$$Neg\_deviat = Max \{ 0, my\ contribution - my\ partner's\ contribution \}$$

$$Pos\_deviat = Max \{ 0, my\ partner's\ contribution - my\ contribution \}$$

The results from the two models are very similar and show that both variables have a positive effect on the probability of punishing the partner. The type of punishment, the type of matching, and the interactive effect of both are also statistically significant. The first two variables have positive sign meaning that the probability of punishing the partner is higher in the STP treatment (in comparison with the SOP) and in the Partner M treatment (in comparison with the Stranger M). We also find that the interaction between STP and Partner M is significant. To calculate the correct interaction effect for this non-linear model we follow Norton et al. (2004). The correct value is -0.2425, and its standard error is 0.097, leading to a z value of -2.49. It implies that the interaction effect is significant and that the effect of being in a situation with strong punishment possibilities is more important in the Stranger M than in the Partner M. The probability of punishing is higher in the STP than the SOP under Stranger, but the difference is almost zero under Partner M. With respect to the deviation variables we find that the greater the negative deviation of the partner's contribution from a subject's own one, the more the subject sanctions. What is perhaps more surprising is that the variable *Pos\_deviat* is statistically significant and has a positive sign. It can be interpreted as evidence of spiteful preferences on the part of some players.<sup>10</sup> However, the effect of a negative deviation is more than 3.5 times as large as that of a positive deviation, and the effect is found to be significant only at the 10% level.

Fehr and Gächter (2000) found also that a subject is more heavily punished the more his contribution falls below the average contribution of other group members —they have groups of 4 subjects. Masclet et al. (2003) found the same pattern of punishment. They also found that subjects that contributed low amounts were using the punishment more number of times than other subjects.

---

<sup>10</sup> According to Falk et al. (2005) "Spiteful sanctions are those that occur because the sanctioning subject values the payoff of the sanctioned subject negatively, regardless of whether the sanctioned subject behaved fairly or not".

We summarize our findings about the strong effects on punishment in the following regularity.

| I_Punished             | Probit estimation |          | Tobit estimation |          |
|------------------------|-------------------|----------|------------------|----------|
|                        | Coef.             | Std.Err. | Coef.            | Std.Err. |
| Strong_P***            | 0.7320            | 0.2326   | 1.2892           | 0.4280   |
| Partner_M**            | 0.6545            | 0.2762   | 1.1647           | 0.5207   |
| Strong_P*Partner_M***  | -0.9659           | 0.3531   | -1.6919          | 0.6379   |
| Neg_deviat***          | 0.5216            | 0.0791   | 0.9016           | 0.1520   |
| Pos_deviat*            | 0.0942            | 0.0534   | 0.1680           | 0.0937   |
| Constant***            | -1.5738           | 0.1932   | -2.6873          | 0.3837   |
| Number of observations | 1880              |          | 1880             |          |
| Wald chi2(5)           | 59.53             |          | 47.83            |          |
| Prob > chi2            | 0.0000            |          | 0.0000           |          |

\*\*\* significant at 1 percent level  
 \*\* significant at 5 percent level  
 \* significant at 10 percent level

Table 7: Estimations for the punishment behavior

*Regularity 2: Probit and logit estimations show that punishment is used more frequently in the STP, especially in the Stranger M treatment, and for negative deviations from the partner’s contribution.*

### 4.3 Final earnings

The subjects were paid —as mentioned before— for the sum of earnings in the 20 periods — converted into euros— plus the show-up fee of 3 euros. In table 8 we show the average final earnings without the show-up fee. The average earnings in the Stranger and Partner M are statistically different<sup>11</sup>, but the difference between STP and SOP is not statistically significant. This difference is, however, significant in the Stranger M treatment, where the average earning is higher in the SOP than in the STP.<sup>12</sup>

|                  |          | Type of Punishment |        | Total |
|------------------|----------|--------------------|--------|-------|
|                  |          | Soft               | Strong |       |
| Type of Matching | Stranger | 10.84              | 9.82   | 10.33 |
|                  | Partner  | 11.76              | 12.25  | 12.00 |
| Total            |          | 11.34              | 11.09  | 11.22 |

Table 8: Average final earnings

If we observe the final earnings in figure 6<sup>13</sup>, we notice that in the Partner M the final earnings are higher in the case when strong punishment was available —except for the last

<sup>11</sup>The p-value of the t-test and the Mann-Whitney U test is 0,000.

<sup>12</sup>The p-value of the t-test and the Mann-Whitney U test is 0,000.

<sup>13</sup>The values in the figure are in tokens, not in euros.

period. But the opposite is observed in the Stranger M treatment, a difference with Fehr and Gächter (2000). They found that punishment opportunities initially cause a relative payoff loss, but towards the end there is a relative payoff gain in both types of matching. Masclet et al. (2003) found that in the Partner treatment both types of sanction increase average earnings in the first five periods, but in the last five both types of punishment generate similar earnings.

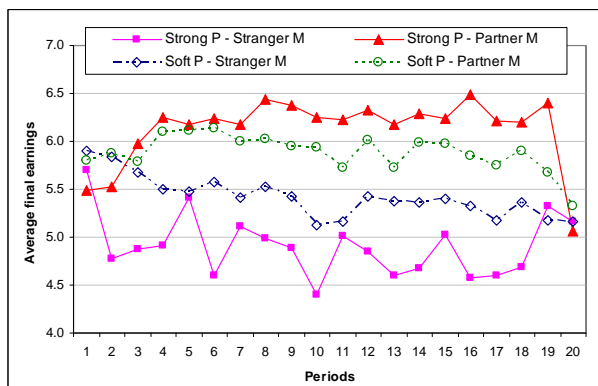


Figure 6: Final earnings in tokens by type of P and M

In the Stranger M the contributions to the group account are almost the same —except for the last periods— in both types of punishment, while the average applied punishment is higher in the STP, resulting in a lower earning in this case in comparison with the SOP. In the Partner M, the contributions are higher in the STP than in the SOP, after period 8, and the average applied punishment is not very different, resulting in higher earnings in the STP.

*Regularity 3: Average final earning are higher in the STP under Partner M, but they are higher in the SOP under Stranger M.*

In figure 7 we have the individual final earnings per period after the punishment is applied, and it shows that the final earnings are very similar for both types of P. The figure on the right shows that the average final earnings are higher in the Partner Matching Treatment than in the Stranger one, as could be expected.<sup>14</sup>

<sup>14</sup> Average Final earnings in period 20 are very similar for all four treatments, ranging between 5.06 and 5.32 tokens.



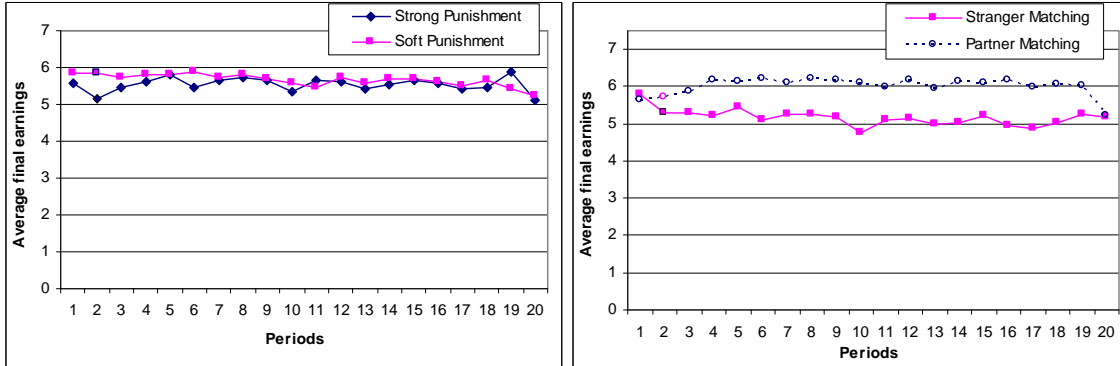


Figure 7: Final earnings by type of P (left) and type of M (right)

## 4.4 Well-being

This section presents the data about the main focus of the paper, namely, the relation between subjective well-being and the type of punishment the subjects are playing under. In section 4.4.1 we show some descriptive statistics about the distribution of the well-being variable and its relation to the types of punishment, matching and to other variables. In section 4.4.2 we present the results of an ordered probit and ordered logit models, in which we estimate the effect of the punishment environment, controlling for relevant variables.

### 4.4.1 Some descriptive statistics

As mentioned above, we measure subjective well-being through self-assessments of participants' satisfaction with the experience in the experiment. The question is: *How satisfied would you say that you are with the experiment?* The possible answers were : 1) Completely satisfied, 2) Very satisfied, 3) Rather satisfied, 4) Neither satisfied nor dissatisfied, 5) Rather dissatisfied, 6) Very dissatisfied, 7) Completely dissatisfied. The aggregate frequencies of answers to the question about well-being (hereafter, WB) can be seen in figure 8. Figure 8 shows that nobody chose level 7 (*Completely dissatisfied*) as an answer, and that people have a tendency to locate in the middle of the distribution. Next we describe the relations between well-being and some other relevant variables.

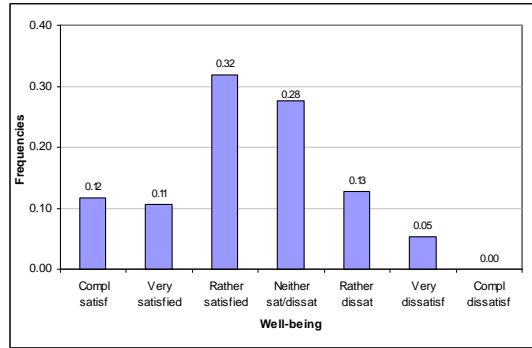


Figure 8: Distribution of the variable WB

**Type of punishment and of matching** The answers to the well-being question by type of punishment are shown in figure 9. This is the main comparison we are interested in.

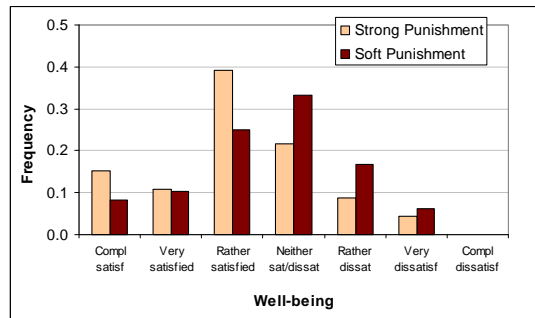


Figure 9: Distribution of the variable WB by type of P

As can be seen in figure 9 the distribution of the variable WB in the STP Treatment is moved to the left with respect to the distribution in the SOP Treatment, indicating that well-being is higher in the first case. The average well-being in SOP is 3.6, and the average in STP is 3.1.<sup>15</sup>

The answers to the well-being question by type of matching are shown in figure 10.

The average well-being in the Stranger M is 3.70, and in the Partner M is 3.04, meaning that the subjects are better, in well-being terms, in the Partner M.<sup>16</sup>

<sup>15</sup>If one takes each observation to be independent, then these averages can be said to be statistically independent. The p-value of the Mann-Whitney U test is 0.061.

<sup>16</sup>The p-value of the t-test is 0.014, and the p-value of the Mann-Whitney U test is 0.000.

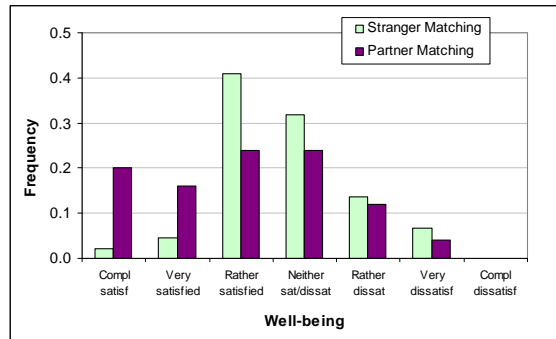


Figure 10: Distribution of the variable WB by type of M

**Other variables** The next step is to ask what **other variables** can explain the WB of subjects.<sup>17</sup> The variable *Final Earnings* in figure 11 has, as could be expected, a positive impact over the well-being of the subject, a higher level of final earnings goes with a lower value of WB and therefore a better well-being.

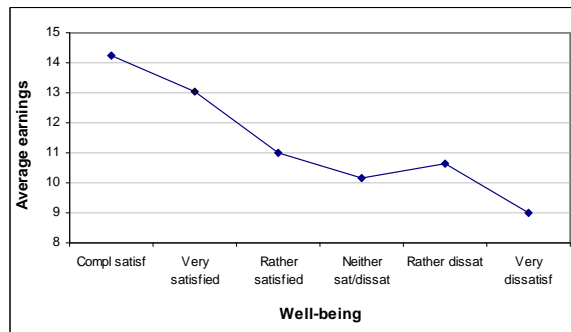


Figure 11: Relation between WB and final earnings

In figure 12 we can see the relation between the different values taken by the variable WB (except for *Very dissatisfied* that nobody chose) and the average values of the variables: average own contribution, average partner contribution and *Me more*. The variable *Me more* captures the number of times a subject contributed more tokens to account A than his partner, regardless of the intensity of the difference.

<sup>17</sup>We gathered some independent information about these other variables. After the experiment we had subjects fill out a questionnaire, where we asked if they agreed or disagreed with some statements, some of them referring to their answer about the satisfaction with the experiment. The results of this questionnaire show that 52.1% of the subjects were influenced by their profit when answering the well-being question, where we counted the subjects that said that “Strongly agree” or “Agree”. 66% said that their partner contribution to account A was influencing their answer, and 31.9% said that the received punishment was influencing.

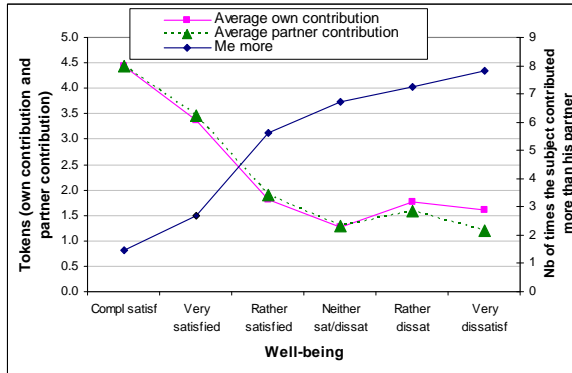


Figure 12: Relation between WB and contributions

The higher the number of times the subject contributed more than his partner, the worse his well-being. On the other hand, the higher his own and his partner contribution, the better in well-being terms.<sup>18</sup> We next describe the relation between WB and punishment received. The variable *Times P* represents the number of times a subject was punished by his partner. We prefer it to the average punishment received because the effective punishment subjects could impose was different depending on the treatment and, in our estimation below, we wanted to estimate one equation including both treatments, moreover, both variables follow a similar pattern and have similar relation with the dependent variable, as can be seen in figure 13.

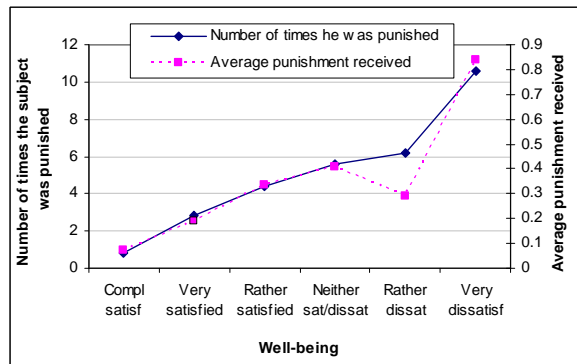


Figure 13: Relation between WB and punishment received

<sup>18</sup>Those with better well-being have higher own contribution and higher partner contribution, that lead to higher earnings.

#### 4.4.2 Estimation of the WB equation

**Econometric model** The variable WB is a discrete variable that can take a value in the 1-7 range, and is ordered from the best WB to the worst. The most commonly used and appropriate methods for estimating models with more than two outcomes, when the dependent variable associated with the outcomes is both discrete and ordinal, are those of ordered logit and ordered probit.

**Our models** The variables included in our estimations for the WB are the following: *Partner M*—a dummy variable that takes value 1 if the observation corresponds to the Partner Matching treatment and 0 otherwise—, *Strong P*—a dummy variable that takes value 1 if the observation corresponds to the Strong Punishment treatment and 0 if it corresponds to the Soft Punishment treatment—, *Final E* that represents the total final earnings of the subjects, *Me\_more* that represents the number of times the subject contributed more tokens to account A than his partner, *Times\_P* that corresponds to the number of times the subject was punished by his partner, and *Times IP*—the number of times the subject punished his partner.

Table 9 shows the results of our estimations. In relation with the goodness-of-fit of the models, it can be said that the null hypothesis that the models did not have greater explanatory power than an "intercept only" model, is rejected.

| WB                     | Ordered probit estimation |          | Ordered logit estimation |          |
|------------------------|---------------------------|----------|--------------------------|----------|
|                        | Coef.                     | Std. Err | Coef.                    | Std. Err |
| Partner M              | 0.2351                    | 0.2687   | 0.5669                   | 0.4900   |
| Strong P ***           | -0.8446                   | 0.2399   | -1.4585                  | 0.4206   |
| Final E ***            | -0.4232                   | 0.1010   | -0.7202                  | 0.1795   |
| Me_more ***            | 0.1068                    | 0.0351   | 0.1986                   | 0.0652   |
| Times_P **             | 0.0789                    | 0.0332   | 0.1439                   | 0.0562   |
| Times IP **            | -0.0583                   | 0.0279   | -0.1095                  | 0.0480   |
| Number of observations | 94                        |          | 94                       |          |
| LR chi2(5)             | 71.40                     |          | 72.48                    |          |
| Prob > chi2            | 0.0000                    |          | 0.0000                   |          |
| Pseudo R2              | 0.2333                    |          | 0.2368                   |          |

\*\*\* significant at 1 percent level

\*\* significant at 5 percent level

Table 9: Results of the oprobit and ologit models

To analyze table 9, it should be taken into account that the dependent variable (*WB*) goes from 1—the most satisfied—to 7—the most dissatisfied—, therefore a positive (negative) sign of the estimated coefficient means that the correspondent variable has a negative (positive) effect over the well-being of the subject.

The two models show very similar results in terms of which variables have a significant impact and with respect to the relative magnitudes of the coefficients. The variable representing the *Type of Matching* is not statistically significant<sup>19</sup>, all other variables are statistically significant and have the expected sign. The variables *Strong P*, *Final Earning*, and *Times IP* have a negative sign, meaning that a higher value of these variables implies a better well-being of the subject, as expected. The variables *Me more* and *Times P* representing the number of times the subject contributed more tokens to account A (group account) and the number of times the subject was punished, respectively, have a negative effect over the well-being. What is crucial here is that controlling for the earnings, contributions and punishment, the variable reflecting the type of punishment is *still* significant. Its negative sign means that in the Strong Punishment situation the subjects have a lower value of the variable WB and therefore a better well-being in comparison with the situation with Soft Punishment.

*Regularity 4: Subjects experience higher subjective well-being under strong than under soft punishment, controlling for other relevant variables.*

In both the oprobit and ologit estimates, the coefficient for strong punishment is about twice as large in magnitude as the one corresponding to the final earnings variable. To confirm whether this first impression is a solid one, we computed the marginal effects from the ordered probit estimates shown in table 9. The marginal effects shown in table 10 reflect increases in probability of being at one of the six well-being levels due to a change in each of the exogenous variables. With respect to the Strong P variable the figures in the table correspond to the effect of the switch from being in the soft punishment environment to being in the strong punishment environment. The table show that switching to the strong punishment environment makes being in each of the three higher well-being levels significantly more likely, while it makes it significantly less likely to be at either of the lower well-being level WB=4 and WB=5.

For the Final E variable, the figures in table 10 correspond to the effect of taking the average individual<sup>20</sup> and increasing that person's earnings by one token. Observe that for each level of WB the marginal effect of the punishment dummy is twice as large as the one for the final

---

<sup>19</sup>We also estimated models including an interactive effect between the type of matching and the type of punishment, and found that the effect is not significant for both models.

<sup>20</sup>In terms of the explanatory variables, i.e. the marginal effects are valued in the mean of the independent variables.

earnings variable. It is with respect to this comparison of the marginal effects that one can make the statement summarized in the following regularity.

|                  | Compl<br>satisfied | Very<br>satisfied | Rather<br>satisfied | Neither<br>sat/dissat | Rather<br>dissatisfied | Very<br>dissatisfied |
|------------------|--------------------|-------------------|---------------------|-----------------------|------------------------|----------------------|
| <b>Variables</b> | <b>P(WB=1)</b>     | <b>P(WB=2)</b>    | <b>P(WB=3)</b>      | <b>P(WB=4)</b>        | <b>P(WB=5)</b>         | <b>P(WB=6)</b>       |
| Partner M (ç)    | -0.0142            | -0.0334           | -0.0421             | 0.0534                | 0.0291                 | 0.0073               |
| Strong P (ç)     | 0.0551*            | 0.1180***         | 0.1415***           | -0.1808***            | -0.1046***             | -0.0292              |
| Final E          | 0.0251**           | 0.0599***         | 0.0773***           | -0.0962***            | -0.0529***             | -0.0133              |
| Me_more          | -0.0063*           | -0.0151**         | -0.0195**           | 0.0243***             | 0.0133***              | 0.0034               |
| Times_P          | -0.0047            | -0.0112**         | -0.0144*            | 0.0179**              | 0.0099**               | 0.0025               |
| Times IP         | 0.0035             | 0.0083*           | 0.0107*             | -0.0133**             | -0.0073*               | -0.0018              |

(ç) dy/dx is for discrete change of dummy variable from 0 to 1

\*\*\* significant at 1 percent level

\*\* significant at 5 percent level

\* significant at 10 percent level

Table 10: Marginal effects of the oprobit model

*Regularity 5: The impact on well-being of the punishment environment is twice as large as that of increasing subjects' earnings.*

Confirmation of the above result can be obtained by analyzing the effects of the dummy variable for the type of punishment in a different way. We do this for the ordered probit model by comparing the estimated probabilities of being at the different WB levels (1, ..., 7) that result when the variable (*Strong P*) takes one value (*Strong P=1*) with the estimated probabilities that are the consequences of it taking the other value (*Strong P=0*), the values of the other variables remaining unchanged between the comparison<sup>21</sup>; using the logit estimates produces very similar results. More specifically, we calculated the probabilities for every subject of being at the different WB levels when *Strong P* was equal to 1 and we calculated the means of the two sets of probability estimates. Formally, if  $\widehat{p}_{ij}^{st}$  and  $\widehat{p}_{ij}^{so}$  are the computed probabilities of person  $i$  being at WB level  $j$  when he is in the *Strong P* (St) and *Soft P* (So) situation respectively, then the means we use are simply<sup>22</sup>

$$\overline{p}_{j}^{st} = \frac{\sum_{i=1}^N \widehat{p}_{ij}^{st}}{N} \quad \overline{p}_{j}^{so} = \frac{\sum_{i=1}^N \widehat{p}_{ij}^{so}}{N}$$

The difference between these probabilities ( $\overline{p}_{j}^{st} - \overline{p}_{j}^{so}$ ) measures the effect of the type of punishment on the mean probability of being at different WB levels. In figure 14 we show the

<sup>21</sup>See Borooah (2002).

<sup>22</sup>To calculate the probabilities we estimated again the oprobit model including only the statistically significant variables. The results are in the appendix.

difference between the probabilities; recall that all the other variables are being kept constant. What is shown in figure 14 confirms the main result of our paper. The Strong Punishment environment makes it more likely that on average people are at the high well-being levels and less likely that they are at the low well-being levels.

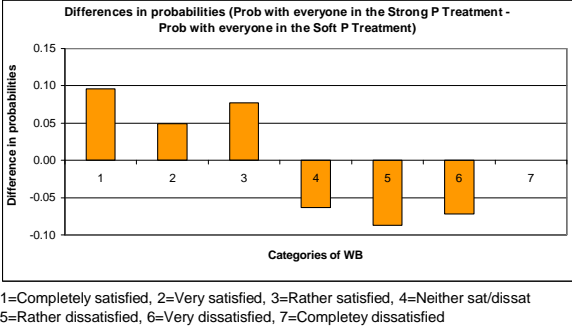


Figure 14: Differences in probabilities (Prob with everyone in the Strong P treatment – Prob with everyone in the Soft P treatment)

## 5 Conclusions

We set out to find a downside to the possibility of using punishment to deter free-riding. Instead, we find that people derive process satisfaction from interacting in a more and not in a less repressive environment. In addition, the positive well-being effect that we find is relatively large. The marginal effect of interacting under strong punishment is about twice as large as the marginal effect of increasing earnings. It turns out that in our environment average earnings do not differ much between the strong and the soft punishment treatments, so that, in our case, the "superiority" of the strong punishment setting stems mostly from process satisfaction considerations. In a context —like in Fehr and Gächter (2000)— where strong punishment possibilities led to higher monetary earnings both effects would go in the same direction.

Our results bear some relation to the study on the neural basis of altruistic punishment by de Quervain et al. (2004), in which subjects’ brains were scanned while they learned about the defector’s abuse of trust and determined the punishment. It was found that a punishment that reduced the defector’s economic payoff activated the dorsal striatum which has been implicated in the processing of rewards that accrue as a result of goal-directed actions. The authors’ interpretation of these results is that people derive satisfaction from punishing norm violations



and that the activation in the dorsal striatum reflects the anticipated satisfaction from punishing defectors.

The results of the work we present in this paper can be seen as independent confirmation of the general notion that people derive satisfaction from punishing. In our work, we study the satisfaction does not derive from the very act of punishing but from being in an environment involving harsh punishment possibilities. Our well-being measure captures the satisfaction derived from the circumstances around the decision-making itself. Note that our result arises in a context in which punishment is sometimes used in a spiteful or at least somewhat unreasonable way; recall that we found that subjects punished their interaction partners for contributing more than themselves. This means that the results we find emerge despite the possibly detrimental effects of such sanctions —see Fehr and Rockenbach (2003).

We believe that our results are of relevance for some very basic issues about the organization of society. In a stable society, social norms guide the interaction among members of society under different circumstances. An important characteristic of organized society is its ability to restrict opportunistic behavior through the use of rewards and punishment. In this way social norms can often support high levels of cooperation among the members of a society.<sup>23</sup> Our results show that, in addition to the material benefits that derive from the possibility of punishing people who do not comply with social norms, people obtain additional satisfaction from interacting under the "protection" of strong punishment possibilities. Güreker et al. (2006) show that a sanctioning institution is the undisputed winner in a vote-with-your-feet type competition with a sanction-free institution. They find that in their experiments the entire population migrates successfully to the sanctioning institution and strongly cooperate, whereas the sanction-free society disappears. Our results suggest that this migration is not only motivated by material payoffs but also by the process-satisfaction of interacting in the sanctioning institution.

A final question is how to find an explanation for the effect we found. Intuition suggests some rather natural ex-post explanations. For example, one could say that people feel better when strong punishment opportunities are available, because they feel "secure"; they feel that the system works. However, to delve deeper into explaining our results one would need additional data. One possibility would be to use some kind of feelings and emotions questionnaire, but perhaps the more fundamental approach to finding an explanation would consist in designing a neural study that would be able to pick up the effect of an environmental variable like the one

---

<sup>23</sup>Dal Bó (2001)

we have studied. Such a study might detect a satisfactory effect of third-party punishment on the brain.

## References

- [1] Blanchflower, David G. and Andrew J. Oswald (2000), "Well-being over time in Britain and the USA". NBER Working paper N°7487. Cambridge, MA: National Bureau of Economic Research.
- [2] Borooah, Vani K.(2002), "Logit and probit. Ordered and multinomial models". Sage University Paper. Series: Quantitative applications in the Social Sciences.
- [3] Brandts, Jordi, Arno Riedl, and Frans van Winden (2004), "Competition and well-being. An Experimental Investigation into Rivalry, Social Disposition, and Subjective Well-Being". Available at <http://www.iae.csic.es/brandts/index.htm>.
- [4] Charness, Gary, and Brit Grosskopf (2001), "Relative Payoffs and Happiness: An Experimental Study," *Journal of Economic Behavior and Organization*, 45, pp. 301-328.
- [5] Dal Bó, Pedro (2001), "Social norms, cooperation and inequality". Working paper N° 802. Department of Economics. University of California, Los Angeles.
- [6] De Quervain, Dominique J.-F., Urs Fischbacher, Valerie Treyer, Melanie Schellhammer, Ulrich Schnyder, Alfred Buck, and Ernst Fehr (2004), "The neural basis of altruistic punishment". *Science*, 27 August 2005, vol. 305, n° 5688, pp. 1254-1258
- [7] Di Tella, Rafael, Robert J. MacCulloch, and Andrew J. Oswald (2001), "Preferences over inflation and unemployment: evidence from surveys of happiness". *American Economic Review*, 91 (1), pp. 335-341.
- [8] Diener, Ed and Shigehiro Oishi (2000), "Money and happiness: income and subjective well-being across nations". In: Diener and Suh (eds.) *Culture and subjective well-being*. Cambridge, MA: MIT Press, pp. 185-218.
- [9] Diener, Ed and Martin E.P. Seligman (2004), "Beyond Money. Toward an Economy of Well-Being". *Psychological science in the public interest*, Volume 5, Number 1.

- [10] Easterlin, Richard (2001), "Income and happiness: toward a unified theory". *Economic Journal*, 111 (473), pp. 465-484.
- [11] Falk, Armin, Ernst Fehr, and Urs Fischbacher (2005), "Driving forces of informal sanctions". *Econometrica*, 73 (6), pp. 2017-2030.
- [12] Fehr, Ernst and Simon Gächter (2000), "Cooperation and punishment in public goods experiments". *American Economic Review*, 90 (4), pp. 980-994.
- [13] Fehr, Ernst and Bettina Rockenbach (2003), "Detrimental Effects of Sanctions on Human Altruism", *Nature*, vol. 422, 13 March 2003, 137-140.
- [14] Fischbacher, Urs (2007), "Z-tree. Zurich toolbox for readymade economic experiments". *Experimental Economics*, 10 (2), pp.171-178.
- [15] Frey, Bruno S., and Alois Stutzer (2000), "What are the sources of happiness?". Working paper N° 0027. Department of Economics. Johannes Kepler University of Linz.
- [16] Frey, Bruno S., and Alois Stutzer (2002), "What can economists learn from happiness research?". *Journal of Economic Literature*, 40 (2). pp. 402-435.
- [17] Gülerk, Özgür, Bernd Irlenbusch, and Bettina Rockenbach (2006), "The Competitive Advantage of Sanctioning Institutions". *Science*, 7 April 2006, vol. 312, n° 5770, 108-111
- [18] Kahneman, Daniel; Ed Diener, and Norbert Schwarz (1999), *Well-being : The foundations of hedonic psychology* Russell Sage Foundation (New York).
- [19] Kahneman, Daniel, Peter P. Wakker, and Rakesh Sarin (1997), "Back to Bentham? Explorations of experienced utility". *The Quarterly Journal of Economics*, 112 (2), pp. 375-405.
- [20] Krueger, Alan B. (2005), "Well-Being and Policy Evaluation". Presentation at the Econometric Society World Congress 2005, London, August 20, 2005.
- [21] McFadden, Daniel (2005), "The New Science of Pleasure". Frisch Lecture, Econometric Society World Congress 2005, London, August 20, 2005.
- [22] Masclet, David, Charles N. Noussair, Steven J. Tucker, and Marie-Claire Villeval (2003), "Monetary and non-monetary punishment in the voluntary contributions mechanism". *American Economic Review*, 93 (1), pp. 366-380.

- [23] Norton, Edward C., Hua Wang, and Chunrong Ai (2004), "Computing interaction effects and standard errors in logit and probit models". *The Stata Journal*, 4 (2,) pp. 154-176.
- [24] Ostrom, Elinor, James Walker, and Roy Gardner (1992), "Covenants with and without a sword: self - governance is possible". *The American Political Science Review*, 86 (2), pp. 404-417.
- [25] Van Praag, Bernard and Ada Ferrer-i-Carbonell (2004), "Happiness Quantified. A Satisfaction Calculus Approach" Oxford University Press.
- [26] Wooldridge, Jeffrey M. (2002), *Econometric analysis of cross-section and panel data*. MIT Press, Cambridge, Massachusetts.

## 6 Appendix 1

The estimates of the ordered probit model for the variable WB including only the statistically significant variables are in table A1.

| Number of observations             | 94      |          |
|------------------------------------|---------|----------|
| LR chi2(5)                         | 71.40   |          |
| Prob > chi2                        | 0.0000  |          |
| Pseudo R2                          | 0.2333  |          |
| WB                                 | Coef.   | Std. Err |
| Strong P ***                       | -0.8374 | 0.2391   |
| Final E ***                        | -0.3886 | 0.0927   |
| Me_more ***                        | 0.1008  | 0.0344   |
| Times_P **                         | 0.0813  | 0.0328   |
| Times IP *                         | -0.0516 | 0.0266   |
| * significant at 10 percent level  |         |          |
| ** significant at 5 percent level  |         |          |
| *** significant at 1 percent level |         |          |

Table A1: Results of the oprobit model

## 7 Appendix 2: Instructions: *Partner Matching, Strong Punishment Treatment*

Thank you for coming to this experiment on decision making. You will be paid 3 euros for showing up plus the money you earn during the experiment which will depend on your and other participants' decisions. At the end of today's session you will be privately paid.

From now on the communication with other participants is not allowed. If you have any doubt during the reading of these instructions or at any moment of the experiment, rise your hand and you will be personally attended.

The experiment consists of two stages.

### *First stage*

The first stage consists of 20 rounds during which you will be randomly paired with another participant, and nobody will know with whom he is playing. During the 20 rounds your partner will be the same participant.

In each round you will be credited with 5 tokens that you will have to decide how to divide (in integers) between two accounts: account A and account B. The tokens allocated to account B will go directly to your earnings, but the tokens allocated to account A will affect your earnings as well as your partner's.

For each token allocated to account A you will receive 0.75 tokens. Your partner will receive the same amount for each token you allocate to account B. In a similar manner, you will receive 0.75 tokens for each token your partner allocates to account A.

The tokens allocated to account B will affect only your earnings, your partner will not receive anything for them, as you will receive nothing for the tokens your partner allocates to his account B.

Therefore, your earnings will be the number of tokens that you allocate to your account B plus the returns from your and your partner's tokens allocated to account A, i.e. Your Earnings = amount of tokens allocated to account B + 0.75 \* total tokens allocated by you and your partner to account A.

The possible earnings are represented in table A2, depending on your contribution (X) and your partner's contribution (Y) to account A —the contributions to account B is 5 minus the contribution to account A. The numbers in the cells represent your earnings (first value) and your partner's earnings (second value). These earnings include the tokens assigned to account B plus the returns from the tokens assigned to account A by you and your partner.

| X\Y | 0           | 1           | 2           | 3           | 4           | 5           |
|-----|-------------|-------------|-------------|-------------|-------------|-------------|
| 0   | 5 , 5       | 5.75 , 4.75 | 6.5 , 4.5   | 7.25 , 4.25 | 8 , 4       | 8.75 , 3.75 |
| 1   | 4.75 , 5.75 | 5.5 , 5.5   | 6.25 , 5.25 | 7 , 5       | 7.75 , 4.75 | 8.5 , 4.5   |
| 2   | 4.5 , 6.5   | 5.25 , 6.25 | 6 , 6       | 6.75 , 5.75 | 7.5 , 5.5   | 8.25 , 5.25 |
| 3   | 4.25 , 7.25 | 5 , 7       | 5.75 , 6.75 | 6.5 , 6.5   | 7.25 , 6.25 | 8 , 6       |
| 4   | 4 , 8       | 4.75 , 7.75 | 5.5 , 7.5   | 6.25 , 7.25 | 7 , 7       | 7.75 , 6.75 |
| 5   | 3.75 , 8.75 | 4.5 , 8.5   | 5.25 , 8.25 | 6 , 8       | 6.75 , 7.75 | 7.5 , 7.5   |

Table A2: Earnings depending on contributions (X\Y) to account A

For example:

- if you allocate 3 tokens to account A and 2 tokens to your account B, and your partner assigns 1 token to account A and 4 tokens to his account B, your earnings will be 5 tokens ( $2 + (3+1)*0.75 = 5$ ) and your partner's earnings will be 7 tokens ( $4 + (3+1)*0.75 = 7$ ). This is represented in the cell corresponding to X=3 and Y=1.
- if you allocate 2 tokens to account A and 3 tokens to your account B, and your partner assigns 2 tokens to account A and 3 tokens to his account B, your earnings will be 6 tokens

$(3 + (2+2)*0.75 = 6)$  and your partner's earnings will be 6 tokens  $(3 + (2+2)*0.75 = 6)$ . This is represented in the cell corresponding to  $X=1$  and  $Y=2$ .

- if you allocate 1 token to account A and 4 tokens to your account B, and your partner assigns 4 tokens to account A and 1 token to his account B, your earnings will be 7.75 tokens  $(4 + (1+4)*0.75 = 7.75)$  and your partner's earnings will be 4.75 tokens  $(1 + (1+4)*0.75 = 4.75)$ . This is represented in the cell corresponding to  $X=1$  and  $Y=4$ .

After every participant has chosen how much to assign to each account, you will be informed about the decision of your partner, and your and your partner's earnings. Then, each participant will have the opportunity of showing his approval/disapproval about his partner's contribution choosing a number of tokens to be deducted from his partner's earnings with certain cost. In table A3 you have the cost (in tokens).

|                   |      |       |       |       |       |       |       |       |       |       |       |       |       |
|-------------------|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Deducted tokens   | 0    | 0.5   | 1.00  | 1.50  | 2.00  | 2.50  | 3.00  | 3.50  | 4.00  | 4.50  | 5.00  | 5.50  | 6.00  |
| Cost of assigning | 0.00 | 0.125 | 0.250 | 0.375 | 0.500 | 0.625 | 0.750 | 0.875 | 1.000 | 1.125 | 1.250 | 1.375 | 1.500 |

Table A3

That means that, for example, to deduce 0.50 tokens from your partner's earnings would cost you 0.125 tokens, to deduce 1 token from your partner's earnings would cost you 0.25 tokens, to deduce 1.50 tokens from your partner's earnings would cost you 0.375 tokens, to deduce 2 tokens from your partner's tokens would cost you 0.50 tokens, to deduce 2.50 tokens would cost you 0.625 tokens, etc.

After this step, the earnings will be again calculated, and they will be equal to:

Final earnings = Initial earnings (from table A2) - Tokens deducted by your partner - Cost of deducing tokens from your partner's earnings (from table A3)

Once the 20 rounds are over we will add up the tokens you have earned in all the rounds and we will calculate the total earnings in euros that will be = total earned tokens \* 0.10, i.e. each token worth 0.10 euros. Therefore, your final payment will be: (Total tokens \* 0.10) euros + 3 euros.

*Second stage*

In this stage a question will be asked. After that, we will pay.