

# Reciprocity: Experiments

## Lecture 18, *Experimental Econ. & Psychology*

Pierpaolo Battigalli  
Bocconi University

12 November 2020

## Abstract

Intention-based reciprocity theory assumes that people wish to be kind towards those they perceive to be kind, and unkind towards those they perceive to be unkind, where kindness depends on intentions, and perceived kindness on the perception of intentions. Therefore, it is a PGT-based theory. Here we consider two experiments, one (Dhaene & Bouckaert 2010) tests the sequential reciprocity model of Dufwenberg and Kirchsteiger (2004), the other (Dufwenberg, Smith & Van Essen 2013) tests the negative reciprocity model in the Hold-Up mini-Game. We start with the latter.

- Many papers experimentally tested other-regarding preferences (see, e.g., the survey by Cooper & Kagel 2016), which include reciprocity as a prominent motivation. Many experiments suggest that models of mere distributional preferences (such as partial altruism or inequity aversion) do not explain well the results because intentions matter.
- Yet, only few papers specifically tested intention-based (hence, belief-dependent) models of reciprocity. Here we focus on:
  - A test of negative (sequential) reciprocity in two versions of the Hold-Up mini-Game by Dufwenberg, Smith & Van Essen (2013).
  - A test of sequential reciprocity theory (Dufwenberg & Kirchsteiger) by Dhaene & Bouckaert (2010).

## Negative reciprocity: DS&V-E

- According to *negative reciprocity theory*, players meet unkindness with unkindness, but (positive) kindness does not matter. Let  $[x]^- = \min\{0, x\}$ ; then, in leader-follower game forms,

$$\begin{aligned}u_1(a_1, a_2, \alpha_{12}) &= \pi_1(a_1, a_2) + \theta_1 \kappa_{12}(a_1, \alpha_{12}) [\kappa_{21}(a_1, a_2)]^-, \\u_2(a_1, a_2, \alpha_{12}) &= \pi_2(a_1, a_2) + \theta_2 [\kappa_{12}(a_1, \alpha_{12})]^- \kappa_{21}(a_1, a_2).\end{aligned}$$

- Dufwenberg, Smith & Van Essen (2013) derive interesting predictions about *hold-up problems* by extending negative reciprocity theory to 3-stage game forms where:
  - pl. 1 can **invest** to produce a good or service at cost  $c$ , or stay **out**; a non-binding contract (e.g., due to unverifiable quality) specifies price  $p > c$ ;
  - pl. 2 can **pay**  $p$ , thus complying with the contract, or renegotiate, holding 1 up with a **take-it-or-leave-it** offer  $t < c$ ;
  - pl. 1 can *accept* (**yes**) or *reject* (**no**).
  - The good/service has no value for pl. 1 and value  $v > p$  for pl. 2.

# Negative reciprocity: experiment of Hold-up mini-Game

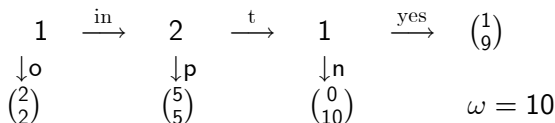
- In the experiment, each player has an endowment (show up fee) of \$2,  $c = 2$ ,  $p = 5$ ,  $t = 1$ ,  $v = 8$ , the resulting game for is as follows:

$$\begin{array}{ccccc} 1 & \xrightarrow{\text{in}} & 2 & \xrightarrow{t} & 1 & \xrightarrow{\text{yes}} & (1) \\ \downarrow o & & \downarrow p & & \downarrow n & & (9) \\ \begin{pmatrix} 2 \\ 2 \end{pmatrix} & & \begin{pmatrix} 5 \\ 5 \end{pmatrix} & & \begin{pmatrix} 0 \\ \omega \end{pmatrix} & & \end{array}$$

E.g., (in, p) yields  $(2 - 2 + 5) = 5$  for pl. 1, and  $(2 + 8 - 5) = 5$  for pl. 2 (different labels are used in the experiment)

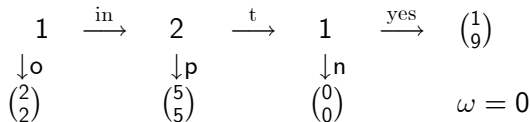
- $\omega = 2 + \text{value for pl. 2 after a rejection}$ :
  - if pl. 1 provided a *service*, he cannot take it back,  $\omega = 2 + 8 = 10$  (**High Game**)
  - if pl. 1 produced a *good* (of no value for him), he can keep it,  $\omega = 2 + 0 = 2$  (**Low Game**)
- According to the *residual right of control*, negative reciprocity yields rejection ( $\omega = 2$ ) promoting cooperation (in, p), or not (if  $\omega = 10$ ).

# Hold-up mini-Game: predictions for the High Game



- **t** is *unkind* after **in**, but rejecting **t** would be a gift of \$1 to pl. 2! Hence, **yes**.
- If pl. 2 anticipates this, he renegotiates with **t** (even if he deems **in** kind, only unkindness is supposed to matter).
- If pl. 1 anticipates this, he stays **out**.
- (Same solution as backward induction with CK of utility=money.)
- Thus, we expect *high rates* of **out**, **t**, **yes**.

# Hold-up mini-Game: predictions for the Low Game



- **t** is *unkind* towards pl. 1 after **in**, rejection hurts pl. 2 a lot ( $-9$ ) and pl. 1 a little ( $-1$ ), for *high enough*  $\theta_1$ , pl. 1's reply is **no**.
- If pl. 2 is afraid of rejection he **complies** (**in** is kind if 1 expects compliance, but only unkindness is supposed to matter).
- If pl. 1 anticipates this, he goes **in**.
- [Under the (preposterous) hypothesis of complete information with high  $\theta$ 's, there is also a "miserable equilibrium" (**in.n, t**) where players are unkind towards each other.]
- Thus, we expect *high rates* of **in, p, no** (the opposite of the High Game).

- Between-subjects experiment (treatments  $\omega = 2$  and  $\omega = 10$ ).
- Subjects were randomly assigned to roles 1 and 2 and played 5 times *in the same role* against changing co-players (hoping to induce some convergence to an equilibrium).
- 5 sessions with 6(H)+6(L) subjects randomly assigned to roles (3 changing pairs in each of H and L) playing 5 rounds:  
 $5 \times 3 \times 5 = 75$  observed plays (terminal histories).
- 1 ECU=1\$.



# Experimental Results (aggr. freq.s in the 75 H/L-plays)

- **High Game**

$$\begin{array}{ccccc} 1 & \xrightarrow{\text{in}} & 2 & \xrightarrow{\text{t}} & 1 & \xrightarrow{\text{yes}} & (1) \\ \frac{45}{75} \downarrow o & & \frac{3}{30} \downarrow p & & \frac{0}{27} \downarrow n & & \omega = 10 \\ \binom{2}{2} & & \binom{5}{5} & & \binom{0}{10} & & \end{array}$$

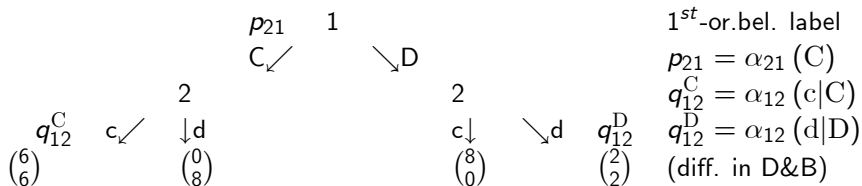
- **Low Game**

$$\begin{array}{ccccc} 1 & \xrightarrow{\text{in}} & 2 & \xrightarrow{\text{t}} & 1 & \xrightarrow{\text{yes}} & (1) \\ \frac{18}{75} \downarrow o & & \frac{20}{57} \downarrow p & & \frac{14}{37} \downarrow n & & \omega = 0 \\ \binom{2}{2} & & \binom{5}{5} & & \binom{0}{10} & & \end{array}$$

- The **null hypothesis** of *treatment-independent* behavior (differences due to randomness) is *rejected* (see pp 9-10 in DS&V-E). The *difference* is in the *predicted direction*.

# Experiment of Dhaene & Bouckaert (2010)

Sequential Prisoners's Dilemma (\$-payoffs)



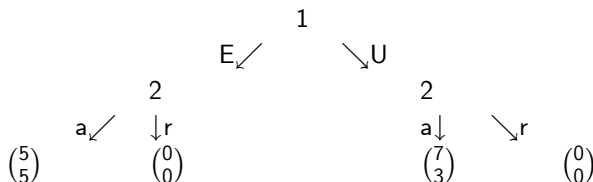
- **PI. 2:** C (resp. D) is certainly kind (resp. unkind), although more kind if  $q_{12}^C$  and  $q_{12}^D$  are low. For sufficiently high  $\theta_2$ , 2's strategy is (c if C, d if D).
- **PI. 1:** let  $\bar{p}_{121} = \mathbb{E}_{\beta_{12}}(p_{21})$  denote pl. 1's (2<sup>nd</sup>-ord.) expectation of  $p_{21}$ ; it can be shown that
  - C is a *material best response* IFF  $6q_{12}^C \geq 8 - 6q_{12}^D$ ;
  - C is a *reciprocity best response* (for suff. high  $\theta_1$ ) IFF  $6\bar{p}_{121}q_{12}^C + 6(1 - \bar{p}_{121})(1 - q_{12}^D) - 3 \geq 0$ .

# Experimental Results of Dhaene & Bouckaert (2010)

- *Look at D&B (2010) pp , 293-294.*
- Average behavior and measured (1<sup>st</sup>- and 2<sup>nd</sup>-order) beliefs: average beliefs are close to unbiased.
- Subjects are classified according to predicted best responses given measured beliefs.
- The behavior of 2-subjects agrees with reciprocity.
- The behavior of 80% of 1-subjects is either consistent with reciprocity, or with selfishness, or both. 20% are “too kind” compared to the reciprocity model.




# Experiment of Dhaene & Bouckaert (2010)

## Ultimatum Game (\$-payoffs)

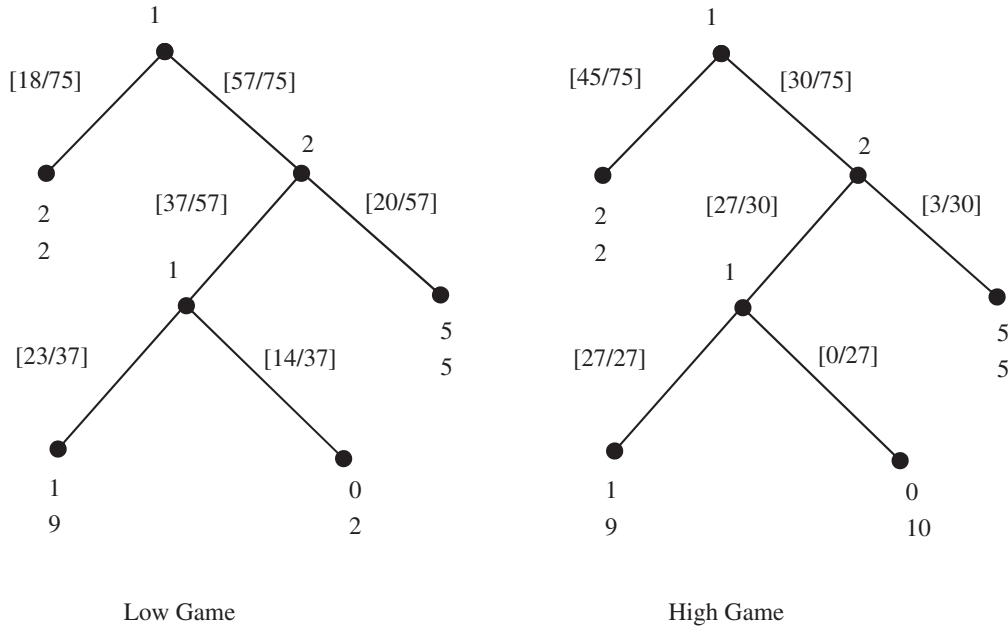


- D&B also analyze the Ultimatum Game shown above. The main difference in results compared to the seq. PD is that average beliefs are biased.

# References

-  BATTIGALLI, P., AND M. DUFWENBERG (2020): “Belief-Dependent Motivations and Psychological Game Theory,” *Journal of Economic Literature*, forthcoming.
-  Cooper, D., and J. Kagel (2016): “Other regarding preferences: A selective survey of experimental results,” *The Handbook of Experimental Economics*, Vol. 2, (J. Kagel and A. Roth eds.) Princeton, PUP.
-  DHAENE, G., AND J. BOUCKAERT (2010): “Sequential Reciprocity in Two-Player, Two-Stage Games: An Experimental Analysis,” *Games and Economic Behavior*, **70**, 289-303.
-  DUFWENBERG, M., AND G. KIRCHSTEIGER (2004): “A Theory of Sequential Reciprocity,” *Games & Economic Behavior*, **47**, 268-298. [Optional.]
-  DUFWENBERG, M., A. SMITH, AND M. VAN ESSEN (2013): “Hold-up: With a Vengeance,” *Economic Inquiry*, **51**, 896-908.

**FIGURE 3**  
Summary of Experimental Results



in the High-treatment. Our research hypothesis at this stage is that *the mean percentage of Y choices is higher in the Low-game than in the High-game*. Table 1 records mean percentage data for the five independent sessions.

A casual look at the data confirms the willingness of subjects to engage in costly punishing once play had reached the third stage. Under the null, the probability of observing a sample as extreme as this one is 0.0027. We therefore clearly reject the associated null. This willingness to punish, even to the detriment of one's own payoff, after player 2 chose action A supports the idea of a negative reciprocity motivation.

The second stage of the game is when player 2 chooses B or A, following player 1's choice of In. Our research hypothesis at this stage, following the discussion in Sections IV.C and V.A, is that *the mean percentage of B choices is higher in the Low-game than in the High-game*.

Table 2 records mean percentage data for the five independent sessions.

Under the null, the probability of observing this sample, or one as extreme, as this one is 0.0362. We thus reject the null. In other words, we find support for the idea that conditional on Player 1 playing In, the efficient equal-split is more likely in the Low-game than the High-game.

At the first stage player 1 chooses whether to trust player 2; In or Out. In the Low-game, player 1 knows he has a punishment mechanism if player 2 chooses A. Our research hypothesis, following the discussion in Sections IV.C and V.A, is that *the mean percentage of In choices is higher in the Low-game than in the High-game*. Table 3 records mean percentage data for the five independent sessions.

Under the null hypothesis that the two samples come from the same distribution, the probability of observing this outcome, or one that is more extreme, is 0.0102. We therefore reject

**TABLE 1**  
Final Stage Choices (Fraction Y)

Treatment	Session 1	Session 2	Session 3	Session 4	Session 5
Low	0.1250	0.4444	0.2500	0.5714	0.6000
High	0.000	0.0000	0.0000	0.0000	0.0000

**TABLE 2**  
Second Stage Choices (Fraction  $B$ )

Treatment	Session 1	Session 2	Session 3	Session 4	Session 5
Low	0.0000	0.2500	0.4286	0.2222	0.6000
High	0.0000	0.1111	0.1429	0.2000	0.0000

**TABLE 3**  
Root Choices (Fraction  $In$ )

Treatment	Session 1	Session 2	Session 3	Session 4	Session 5
Low	0.5333	0.8000	0.9333	0.6000	0.9333
High	0.2667	0.6000	0.4667	0.3333	0.3333

the null. It is plain that in all five sessions the mean percentage of  $In$  choices was higher in the Low-treatment than in the High-treatment.

We noted in Section III that to the extent that the miserable VE described there would have been relevant to the Low-game, negative reciprocity could have been an important motivational force even if there would not have been much of a difference in the nature of play between the High-game and the Low-game. In light of the data, this point now seems moot. All in all, we take the support for our research hypotheses as reinforcing the idea that negative reciprocity can mitigate hold-up mainly in cases where the investing party maintains the residual rights of control.

## VI. CONCLUDING REMARKS

The back cover of the *JPE* once recalled a hold-up story about a rich woman in Savannah where, between the lines, we see negative reciprocity at work<sup>9</sup>:

Some years ago she ordered a pair of iron gates for her house. They were designed and built especially for her. But when they were delivered she pitched a fit, said they were horrible, said they were filth. "Take them away," she said, "I never want to see them again!" Then she tore up the bill, which was for \$1,400—a fair amount of money in those days.

The foundry took the gates back, but didn't know what to do with them . . . there wasn't much demand for a pair of ornamental gates exactly that size. The only thing they could do was sell the iron for its scrap value. So they cut the price from \$1,400 to \$190. Naturally, the following day the woman sent a

man over to the foundry with \$190, and today those gates are hanging on her gateposts where they were originally designed to go.

The story may seem puzzling. Why would the woman send a man to the foundry rather than just make a take-it-or-leave-it offer herself? Part of the answer may be that she feared a counter-offer, but another part is that she might otherwise irritate the foundry's owner who may retaliate by refusing to sell her the gate. On this interpretation, we thus have a situation where a proper understanding of an economic outcome involves reference to negative reciprocity. And if we modify the situation to make the foundry less naive, that is, so that they could see through the woman's ploy, the situation would structurally resemble our Example 1, or our Low-game.

Classical hold-up theory typically assumes that the involved parties selfishly maximize own income. We have argued that this perspective may be too limited; negative reciprocity may plausibly play a role too. Injured parties may have an inclination to strike back if they are treated badly (even if this is costly), and if this is anticipated the problems because of hold-up are mitigated. We have shown, however, that it would be premature to draw the blanket conclusion that hold-up is not a serious concern. Rather, this depends in predictable ways on details of the situation. Namely, hold-up is a less serious concern if the investing party retains residual rights of control than if the other party does. This conclusion is supported by a D&K-based theory of negative reciprocity which we apply to two hold-up games (which vary the residual right of control), and through a related experimental test.

9. See *Journal of Political Economy* 107(1), February 1999. The excerpt is from John Berendt's 1994 novel *Midnight in the Garden of Good and Evil*. It was suggested to the *JPE* by Oliver Hart, and to us by Tore Ellingsen.

**Table 1**

Average behavior and beliefs in the SPD.

A's choice		B's choice	
A's C-rate	0.41 (27/66)	B's c-rate following C	0.37 (10/27)
B's average $p'$	0.35 ( $n = 66$ )	A's average $q'_{c C}$	0.28 ( $n = 66$ )
A's average $p''$	0.44 ( $n = 66$ )	B's average $q'_{c C}$	0.30 ( $n = 66$ )
		B's d-rate following D	1 (39/39)
		A's average $q'_{d D}$	0.84 ( $n = 66$ )
		B's average $q'_{d D}$	0.88 ( $n = 66$ )

participants equally divided between two rooms: room A, where the subjects assumed the role of player A; and room B, where the subjects assumed the role of player B. The subjects played the game once; there was no repeated play or role reversal. The experiment was carried out sequentially, the first part in room A and the second part in room B. This enabled us to elicit B's choice using the direct-response method. That is, the B's responded by choosing  $c$  or  $d$  after observing A's choice,  $C$  or  $D$ . The experiment took about 30 minutes in each room. The A's could not communicate with the B's between the two parts. The subjects were also asked to report their prior beliefs. There were two short questionnaires (one for the A's, one for the B's) with three questions each, measuring

$$(q'_{c|C}, q'_{d|D}, p'') \quad \text{for the subjects in room A,}$$

$$(p', q''_{c|C}, q''_{d|D}) \quad \text{for the subjects in room B.}$$

For example, to measure  $q'_{c|C}$  we asked the subjects in room A [italics added]: "What percentage of people in room B who learned that the person from room A with whom he/she is paired chose option A1 [meaning C] will subsequently choose option B1 [meaning c]?" and to measure  $q''_{c|C}$  we asked the subjects in room B: "What is the average answer of the people in room A to question a1 above [referring to the former question]?" The questions had to be answered after the instructions were given but before playing. A's choice ( $C$  or  $D$ ) was disclosed to the corresponding B after B had answered the questions. To elicit beliefs, a bonus of 3 euros was given for each answer that deviated no more than 5 percentage points from the true percentage. The experiment was set up in such a way that the material payoffs, the questionnaires to both players, the bonus system, the direct-response feature, and anonymity were common knowledge.<sup>11,12,13</sup>

### 2.1.5. Experimental results

The average earnings per subject (including bonuses) were 4.5 euro for player A and 4.92 euro for player B. Table 1 gives the average behavior and beliefs.<sup>14</sup> There is a striking agreement between the average behavior, the average beliefs about behavior, and the average beliefs about those beliefs. This will be analyzed formally in Section 3, where equilibrium behavior is investigated.

Using Proposition 1, the subjects were classified according to their material best response and their reciprocity best response as implied by their beliefs. Table 2 reports the behavioral rates within each category (A's behavior in the upper part; B's behavior in the lower part). For example, the first row shows that 6 A's had beliefs that implied  $C$  was both a material best response and a reciprocity best response; 4 of these A's chose  $C$  and 2 chose  $D$ . These 2  $D$ -choices are neither a material best response nor a reciprocity best response, which is indicated by a "\*". In the group where  $C$  was a material best response and  $D$  a reciprocity best response, 6 out of the 7 subjects chose  $C$ . For these subjects

<sup>11</sup> The questionnaires used to elicit beliefs, inevitably, became part of the framing and hence may have affected behavior. Potential effects of belief elicitation could not be avoided in our set-up, given that B's choice was elicited using the direct-response method, which rules out measuring prior beliefs post-play. There is, however, a potentially important advantage of measuring beliefs prior to play. Belief elicitation, whether pre- or post-play, explicitly invites subjects to reflect on the strategic situation. In case of post-play belief elicitation, subjects may—on second thought, triggered by the questions—change their views of the situation and report beliefs that differ from the (perhaps more vague) beliefs on which their strategic choice was based. Pre-play belief elicitation, using questions as an integral part of the instructions, is more likely to succeed in accurately measuring true beliefs associated with the strategic choices made.

<sup>12</sup> As Charness and Dufwenberg (2006) note, reported beliefs elicited as above may deviate somewhat from true beliefs, defined as the mean of the subject's prior distribution. Thus, for example, rational players would never report beliefs less than 5% or greater than 95%. This sort of measurement error is presumably small. An alternative would have been to use a quadratic scoring rule, which is incentive-compatible with truth-telling if subjects are risk-neutral, but more complicated and subject to Harrison's (1989) flat-maximum critique. Another complication arises from the increased likelihood that subjects make unintentional errors when forming beliefs, compared to the relatively more simple task of choosing between two strategies (e.g.,  $C$  and  $D$ ). For example, it suffices that some A's misinterpret the question asking to report their belief  $q'_{d|D}$  as asking to report  $q'_{c|D} = 1 - q'_{d|D}$  (a careful inspection of the raw data shows that this almost certainly happened) to have a dramatic effect on the average of the reported beliefs. The likelihood of making errors, no doubt, further increases when subjects have to form beliefs about other subjects' beliefs.

<sup>13</sup> Paying subjects for accurate beliefs creates hedging opportunities, i.e. subjects may use reported beliefs to hedge against unfavorable outcomes, potentially biasing reported beliefs and affecting behavior. However, such effects appear to be small. See Blanco et al. (2008), who also propose a hedging-proof belief elicitation method.

<sup>14</sup> The raw data on behavior and beliefs and the material and reciprocity best responses implied by beliefs are available as supplementary material.



**Table 2**  
Best-response analysis of behavior in the SPD.

	A's best response			A's behavior	
	Material	Reciprocity		A's C-rate	A's D-rate
	C	C	(n = 6)	0.67 (4/6)	0.33 (2/6)*
	C	D	(n = 7)	0.86 (6/7)	0.14 (1/7)
	D	C	(n = 6)	0.83 (5/6)	0.17 (1/6)
	D	D	(n = 47)	0.26 (12/47)*	0.74 (35/47)
A's behavior	B's best response			B's behavior	
	Material	Reciprocity		B's c-rate	B's d-rate
C	d	c	(n = 27)	0.37 (10/27)	0.63 (17/27)
D	d	d	(n = 39)	0 (0/39)*	1 (39/39)

\* Indicates behavior that is neither a material best response nor a reciprocity best response.

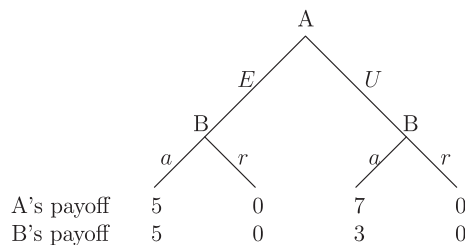
$\kappa_A > 0$  and  $\lambda_A < 0$  (recall that  $\kappa_A$  is A's kindness to B and  $\lambda_A$  is A's belief about B's kindness to A), which, in the context of DK's model, indicates low or no sensitivity to reciprocity. Conversely, when D was a material best response and C a reciprocity best response, 5 out of the 6 subjects chose C. This and the previous observation, taken together, suggest different levels of reciprocity in these two groups of subjects and, since beliefs determine group membership, reciprocity sensitivities appear to be related to beliefs. In the remaining group, where D was both a material best response and a reciprocity best response, 35 out of the 47 subjects chose D (with  $\kappa_A$  and  $\lambda_A$  both negative), and 12 out of the 47 subjects chose C, which is neither a material best response nor a reciprocity best response. For the latter subjects, arguably the most interesting, the typical belief pattern is as follows:  $q'_{c|C}$  is small,  $q'_{d|D}$  is large, and  $p''$  is moderate to large (the medians are 0.18, 0.95, and 0.72). Hence, A believes that B is unkind and that D will yield the highest material payoff, but she still chooses C. The typical beliefs do not suggest obvious unintentional errors, which makes it difficult to understand why these subjects chose C.<sup>15</sup> The results for the A's may be summarized as follows: A's behavior was a material or reciprocity best response in 52 out of the 66 cases (or 79%). Where the predictions of DK's theory failed, the observed behavior was nearly always too kind. That is, negative reciprocation in anticipation was observed less frequently than predicted.

Now consider the B's. Following A's choice of C, B's choice must be either a material best response (d, which occurred 17 times out of the 27) or a reciprocity best response (c, which occurred 10 times out of the 27). Thus, because DK's theory is in line with either of B's choices following A's choice of C, it is non-falsifiable in this particular instance. Following A's choice of D, however, the only possible best response of B is to choose d, thus maximizing her material and reciprocity payoffs. In line with the predictions of DK's theory, this occurred in all of the 39 cases.

## 2.2. The mini-ultimatum game

### 2.2.1. The game

Consider the mini-ultimatum game depicted in Fig. 2. Player A proposes dividing an amount equally (E) or unequally (U); player B observes A's choice and then decides to accept (a) or to reject (r) the proposal. The material payoffs are given at the end nodes.



**Fig. 2.** The mini-ultimatum game.

<sup>15</sup> Maybe other motivations made some subjects choose C, perhaps guilt-aversion or the willingness to give B the benefit of the doubt even when  $q'_{c|C}$  is small. Low values of  $q'_{c|C}$  combined with a C-choice might also suggest unconditional altruism. But then one would equally expect some B's to be unconditional altruists and, when confronted with A's choice of D, to turn the other cheek and choose c; but not a single subject, out of the 39, did so. Furthermore, note the following: An unconditionally altruistic B may choose c following A's choice of C because she is willing to give up 2 to increase A's material payoff by 6; 10 out of the 27 B's did so. But now, following A's choice of D, an unconditionally altruistic B faces a similar question: Will she give up 2 to increase A's material payoff by 6? Nobody did. Thus, we are led to conclude that either unconditional altruism vanishes when both agents exert control over the course of events, or that unconditionally altruistic preferences are highly convex even when comparing very small monetary payoffs pairs.