

Lecture 5: Altruism, Fairness, and Social Preferences

EC 404: Behavioral Economics
Professor: Ben Bushong

April 9, 2024

The Standard Model: Pure Self Interest

The standard model: pure self interest.

- ▶ A person cares only about his or her own outcomes, and not at all about the outcomes of others.

A famous passage from Adam Smith (1776):

It is not from the benevolence of the butcher, the brewer, or the baker that we expect our dinner, but from their regard for their own interest. We address ourselves not to their humanity, but to their self-love, and never talk to them of our necessities, but of their advantage.

A Motivating Quote

A not-as-famous passage from Adam Smith (1759):

No matter how selfish you think man is, it is obvious that there are some principles in his nature that give him an interest in the welfare of others, and make their happiness necessary to him, even if he gets nothing from it but the pleasure of seeing it.

Years before writing *The Wealth of Nations*, Adam Smith had concerned himself with these ideas. Then they lay dormant for a few generations.

A Motivating Anecdote

A passage from Robin Dawes and Richard Thaler (1988):

In the rural areas around Ithaca it is common for farmers to put some fresh produce on the table by the road. There is a cash box on the table, and customers are expected to put money in the box in return for the vegetables they take. The box has just a small slit, so money can only be put in, not taken out. Also, the box is attached to the table, so no one can (easily) make off with the money. We think that the farmers have just about the right model of human nature. They feel that enough people will volunteer to pay for the fresh corn to make it worthwhile to put it out there. The farmers also know that if it were easy enough to take the money, someone would do so.

An Overview of the Literature

Note: The literature on alternatives to pure self-interest is almost entirely an experimentally based literature.

- ▶ Begin with experiments that contradict pure self-interest.
- ▶ Use these experiments to motivate alternative theories.
- ▶ Develop new experiments to test these alternative theories.
- ▶ Use these new experiments to modify alternative theories.
- ▶ And so on....

Experimental Evidence that Contradicts Pure Self Interest

We'll focus on evidence from three experimental games:

- ▶ Prisoners' Dilemma Games
- ▶ Dictator Games
- ▶ Ultimatum Games

For all studies:

- ▶ anonymous interactions
- ▶ “one-shot” interactions
- ▶ typically for real money

Prisoners' Dilemma Games

An example of a Prisoners' Dilemma Game:

	C	D
C	3, 3	0, 4
D	4, 0	1, 1

Key feature: (C, C) maximizes social surplus, but D is a dominant strategy for each player.

Prisoners' Dilemma Games

In Prisoners' Dilemma Games:

- ▶ Pure self interest implies that people should choose D .
- ▶ In experiments, it is not uncommon for people to choose C .
- ▶ For instance, in Cooper *et al* (*GEB* 1996), cooperation rates are about 30%.

Note two other features of data:

- ▶ There is a great deal of heterogeneity — some people cooperate a lot, others cooperate very little.
- ▶ Subjects play 20 rounds with new opponents, and as time passes and they get feedback on earlier games, people cooperate less.

Dictator Games

Dictator Games:

- ▶ There is a pie of size P .
- ▶ Player 1 (the dictator) offers a share $s \in [0, 1]$.
- ▶ Payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$.

In Dictator Games:

- ▶ Pure self interest implies that Player 1 should choose $s = 0$.
- ▶ In experiments, it is not uncommon for people to choose $s > 0$.
- ▶ For instance, in Forsythe *et al* (*GEB* 1994):
 - ▶ about 20% chose $s = 0$
 - ▶ about 60% chose an s between 0 and 0.5
 - ▶ about 20% chose $s = 0.5$

Ultimatum Games

Ultimatum Games:

- ▶ There is a pie of size P .
- ▶ Player 1 (the proposer) offers a share $s \in [0, 1]$.
- ▶ Player 2 (the responder) then accepts or rejects this offer.
- ▶ If Player 2 accepts, payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$; otherwise $x_1 = x_2 = 0$.

Ultimatum Games

In Ultimatum Games:

- ▶ Pure self interest implies that Player 2 should accept any $s > 0$, and so Player 1 should choose the smallest positive s (and perhaps $s = 0$).
- ▶ In experiments, it is not uncommon for Player 1's to choose $s \gg 0$ and for Player 2's to reject low offers.
- ▶ For instance, Fehr and Schmidt (*QJE* 1999) survey 11 prior experiments, and conclude that:
 - ▶ (i) Player 1's virtually never choose $s > 0.5$;
 - ▶ (ii) The vast majority of Player 1's choose an s between 0.4 and 0.5;
 - ▶ (iii) Player 1's rarely choose $s < 0.2$; and
 - ▶ (iv) Player 2's frequently reject low offers, and the probability of rejection is larger the smaller is s .

A Framework for Social Preferences

Let $\mathbf{x} = (x_1, x_2)$ be a vector of material payoffs.

- ▶ x_1 is player 1's material payoff.
- ▶ x_2 is player 2's material payoff.

Pure self interest: Player 1's preferences are $u^1(x_1)$.

Social preferences: Player 1's preferences are $u^1(x_1, x_2)$.

[illegible]

More generally, let $\mathbf{x} = (x_1, x_2, \dots, x_N)$ be a vector of material payoffs.

- ▶ x_i is player i 's material payoff.

Pure self interest: Player 1's preferences are $u^1(x_1)$.

Social preferences: Player 1's preferences are $u^1(x_1, x_2, \dots, x_N)$.

Can Simple Altruism Explain this Evidence?

Definition: A person exhibits *simple altruism* if her utility is increasing in other people's material payoffs (i.e., she likes other people's material payoffs to increase).

A simple formulation:

$$u^1(x_1, x_2) = x_1 + \phi x_2 \quad \text{for some } \phi > 0.$$

\implies trade off own material payoff vs. others' material payoffs.

Altruism in the Prisoners' Dilemma

Reminder: The Prisoners' Dilemma is:

	C	D
C	3, 3	0, 4
D	4, 0	1, 1

Because choosing C rather than D increases the other player's material payoff, simple altruism militates in favor of a person choosing C .

Hence, if a person's degree of altruism is large enough (if ϕ is large enough), it will be optimal to choose C .

Altruism in the Dictator Game

Reminder: The Dictator Game is:

- ▶ There is a pie of size P .
- ▶ Player 1 (the dictator) offers a share $s \in [0, 1]$.
- ▶ Payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$.

Because choosing $s > 0$ increases the other player's material payoff, simple altruism militates in favor of choosing $s > 0$ (more precisely, simple altruism argues for $s = 1$).

Hence, if a person's degree of altruism is large enough (if ϕ is large enough), it will be optimal to choose $s > 0$.

Altruism in the Ultimatum Game

Reminder: The Ultimatum Game is:

- ▶ There is a pie of size P .
- ▶ Player 1 (the proposer) offers a share $s \in [0, 1]$.
- ▶ Player 2 (the responder) then accepts or rejects this offer.
- ▶ If Player 2 accepts, payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$; otherwise $x_1 = x_2 = 0$.

For Player 1, the logic is much as in the Dictator Game.

But for Player 2, rejecting an offer $s > 0$ decreases both the person's own material payoff and the other player's material payoff, and so material-payoff concerns and simple altruism both militate in favor of accepting.

Hence, altruism cannot explain responder rejections in the ultimatum game.

Summarizing our analysis of simple altruism:

- ▶ Although simple altruism can explain some experimental results, it cannot explain everything. Hence, we need something more sophisticated.

“Inequity Aversion” / “Difference Aversion”

Consider the following idea (proposed independently by Fehr & Schmidt (1999) and Bolton & Ockenfels (2000)):

Definition: A person exhibits *inequity aversion* (aka *difference aversion*) if her utility is decreasing in material-payoff differences (i.e., she likes reductions in material-payoff differences).

A simple model (Fehr-Schmidt formulation):

$$u^1(x_1, x_2) = \begin{cases} x_1 - \alpha(x_2 - x_1) & \text{if } x_1 \leq x_2 \\ x_1 - \beta(x_1 - x_2) & \text{if } x_1 \geq x_2 \end{cases}$$

where $0 \leq \beta < 1$ and $\alpha \geq \beta$.

\Rightarrow trade off own material payoff vs. inequity reduction.

Inequity Aversion in the Prisoners' Dilemma

Reminder: The Prisoners' Dilemma is:

	C	D
C	3, 3	0, 4
D	4, 0	1, 1

If you believe the other player will play C , then choosing C rather than D reduces material-payoff differences, and therefore inequity aversion militates in favor of choosing C .

Hence, if a person's degree of inequity aversion is large enough (if β is large enough), it will be optimal to choose C — in which case (C, C) is an equilibrium.

Inequity Aversion in the Dictator Game

Reminder: The Dictator Game is:

- ▶ There is a pie of size P .
- ▶ Player 1 (the dictator) offers a share $s \in [0, 1]$.
- ▶ Payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$.

Because choosing $s > 0$ reduces material-payoff differences, inequity aversion militates in favor of choosing $s > 0$ (more precisely, inequity aversion argues for $s = \frac{1}{2}$).

Hence, if a person's degree of inequity aversion is large enough (if β is large enough), it will be optimal to choose $s > 0$.

Inequity Aversion in the Ultimatum Game

Reminder: The Ultimatum Game is:

- ▶ There is a pie of size P .
- ▶ Player 1 (the proposer) offers a share $s \in [0, 1]$.
- ▶ Player 2 (the responder) then accepts or rejects this offer.
- ▶ If Player 2 accepts, payoffs are $x_1 = (1 - s)P$ and $x_2 = sP$; otherwise $x_1 = x_2 = 0$.

For Player 1, the logic is much as in the Dictator Game.

For Player 2, rejecting an offer $s < 1/2$ reduces material-payoff differences, and thus inequity aversion militates in favor of rejection.

Hence, if a person's degree of inequity aversion is large enough (if α is large enough), it will be optimal to reject small offers.

Inequity Aversion: Conclusion

Summarizing our analysis of inequity aversion:

- ▶ Inequity aversion can indeed explain experimental results in the Prisoners' Dilemma, Dictator Games, and Ultimatum Games (and several other commonly studied games).

A Problem with Inequity Aversion

Are people as destructive as inequity aversion suggests?

More recent research suggests not.

Some simple dictator-game experiments from Charness & Rabin (2002):

(1) Choose between

(400, 400)

30 %

(400, 750)

70%

$[N = 26 \text{ (Berkeley)}]$

(2) Choose between

(400, 400)

50 %

 $(375, 750)$

50%

$[N = 80 \text{ (both locations)}]$

Some Evidence from Andreoni & Miller (2002)

Andreoni and Miller (*ECTA* 2002) studied modified dictator games:

- ▶ Player 1 asked to divide N tokens between Player 1 and Player 2.
- ▶ Each token is worth t_1 points to Player 1, and t_2 points to Player 2.

Main focus: How does the division depend on t_1 and t_2 ?

Some Evidence from Andreoni & Miller (2002)

Three relevant categories of preferences:

(i) Pure self interest: maximize own payoff.

- ▶ In general, $u^1(x_1, x_2, \dots, x_N) = x_1$
- ▶ Here, keep all tokens for yourself.

(ii) Maximin: maximize payoff of least-advantaged person.

- ▶ In general, $u^1(x_1, x_2, \dots, x_N) = \min\{x_1, x_2, \dots, x_N\}$
- ▶ Here, divide tokens in a way that equalizes points.

(iii) Utilitarian: maximize the sum of payoffs.

- ▶ In general, $u^1(x_1, x_2, \dots, x_N) = x_1 + x_2 + \dots + x_N$
- ▶ Here, give all tokens to whomever values them the most.

Some Evidence from Andreoni & Miller (2002)

Results:

Subjects can be roughly classified into these categories:

Table 3 — Classification of Subjects (176 students)

<u>Utility Function</u>	<u>Strong Fit</u>	<u>Weak Fit</u>
Pure Self Interest	40	43
Maximin	25	29
Utilitarian	11	28

Another Model: Social-Welfare Preferences

Adapted from Charness & Rabin (2002)

Definition: A person exhibits *social-welfare preferences* if her utility is increasing in a social-welfare function (i.e., she likes higher “social welfare”).

A simple formulation:

$$u^1(x_1, x_2) = x_1 + \lambda * W(x_1, x_2) \quad \text{for some } \lambda > 0.$$

where

$$W(x_1, x_2) = \delta * \min\{x_1, x_2\} + (1 - \delta) * [x_1 + x_2]$$

- ▶ $\delta = 0 \implies u^1(x_1, x_2) = x_1 + \lambda * (x_1 + x_2)$
- ▶ $\delta = 1 \implies u^1(x_1, x_2) = x_1 + \lambda * \min\{x_1, x_2\}$
- ▶ $\delta \in (0, 1) \implies$ something “in between”.

\implies trade off own material payoff vs. social-welfare concern.

A Problem with Social-Welfare Preferences

We've noted two types of social preferences: constructive (all three on previous screen) and destructive (inequity aversion).

With constructive social concerns, how can we explain rejections in the Ultimatum Game?

An answer: “Intentions” — how we got there matters.

- ▶ Up to now, we have considered only “payoff-based approaches” wherein all that matters are the final outcomes.
- ▶ But perhaps people care about how those final outcomes are reached.

A Motivating Story

A motivating story from Robert Frank (1994):

There is an often-told story of a boy who found two ripe apples as he was walking home from school with a friend. He kept the larger one for himself, and gave the smaller one to his friend. "It wasn't fair to keep the larger one for yourself", the friend replied. "What would you have done?" the first boy asked. "I'd have given you the larger one and kept the smaller one for myself," said the friend. To which the first boy responded, "Well, we each got what you wanted, so what are you complaining about?"

Some Motivating Hypothetical Examples

Some Evidence from Blount (OBHDP 1995)

Consider two versions of the “Ultimatum Game”:

Version 1:

- ▶ Player 1 offers a share

$$s \in \{\$0.00, \$0.50, \$1.00, \dots, \$9.50, \$10.00\}$$

- ▶ Player 2 then accepts or rejects this offer.
- ▶ If Player 2 accepts, payoffs are $x_1 = 10 - s$ and $x_2 = s$; otherwise $x_1 = x_2 = 0$.

Some Evidence from Blount (OBHDP 1995)

Consider two versions of the “Ultimatum Game”:

Version 2:

- ▶ A computer randomly chooses a share

$$s \in \{\$0.00, \$0.50, \$1.00, \dots, \$9.50, \$10.00\}$$

- ▶ Player 2 then accepts or rejects this offer.
- ▶ If Player 2 accepts, payoffs are $x_1 = 10 - s$ and $x_2 = s$; otherwise $x_1 = x_2 = 0$.

Some Evidence from Blount (OBHDP 1995)

Results for Version 1 (where Player 1 chooses the share):

- ▶ Only 29% accept (\$9.50, \$0.50).
- ▶ Only 41% accept (\$8.00, \$2.00).
- ▶ Only 64% accept (\$6.00, \$4.00).

Results for Version 2 (where a computer chooses the share):

- ▶ 80% accept (\$9.50, \$0.50).
- ▶ The other 20% reject (\$6.00, \$4.00).

...that means we still lack a good model for their behavior. Maybe they don't care?

Exactly how do intentions matter?

An often discussed motive is “reciprocity”:

- *Reciprocity* (aka *reciprocal altruism*): You are motivated to be kind to those who are kind to you, and you are motivated to be unkind to those who are unkind to you (and you trade off these motivations against material payoffs).

For instance, in the ultimatum game, making a low offer is an unkind act, and so if I'm a reciprocal altruist I'll be motivated to be unkind back, and hence I might reject.

Charness & Rabin (2002) investigate further the nature of reciprocity....

Evidence of Negative Reciprocity?

Evidence of Negative Reciprocity?

Hence, they do NOT find evidence of “strong reciprocity” wherein a person has a desire to punish someone who behaves unkind — because people are not choosing to take advantage of a free punishment.

Instead, they see evidence of what they label “concern withdrawal”:

- ▶ **Concern withdrawal**: If another player has behaved too selfishly, you remove your concern for that person from your social-welfare preferences.

Evidence of Positive Reciprocity?

Evidence of Positive Reciprocity?

They do NOT find much evidence of positive reciprocity — only when it is free to do so.

Possible interpretation?

- ▶ People ought to pursue good social outcomes, and thus there is no reason to reward them when they do so (because they are just doing what they ought to do).

Two Further Issues

An Issue: Intentions vs. Types vs. Expectations

- ▶ Intentions: Is the other person being nice or nasty to me?
- ▶ Types: Is the other person a decent type or a nasty type?
- ▶ Expectations: What's a reasonable payoff to expect in this game?

Another issue: Social preferences can be highly “context” dependent, because the “context” can influence what is viewed as “fair” or as a “good social outcome”.

- ▶ Consider two examples....

Some Evidence from Krupke & Weber (JEEA 2013)

Consider two versions of a “Dictator Game”:

- ▶ Standard Dictator Game: Player 1 is given \$10. Then Player 1 is given the opportunity to share some with Player 2 (in \$1 increments).
- ▶ Bully Dictator Game: Player 1 and Player 2 are each given \$5. Then Player 1 is given the opportunity to either (i) share some of her \$5 with Player 2 or (ii) take some of Player 2's \$5 (in \$1 increments).

Note: In terms of final payoffs, these games are identical — they both involve Player 1 choosing how to allocate \$10 between Player 1 and Player 2 (in \$1 increments).

Some Evidence from Lazeer, Malmendier & Weber (2012)

Consider two slightly different versions of a “Dictator Game”:

- ▶ Standard Dictator Game: Player 1 is given \$10. Then Player 1 is given the opportunity to share some with Player 2 (in \$1 increments).
- ▶ Dictator Game with Sorting: First, Player 1 decides whether to (i) get \$10 vs. (ii) play a split-\$10 dictator game. If Player 1 chooses option (i), the game ends. If Player 1 chooses option (ii), then Player 1 is given \$10 and given the opportunity to share some with Player 2 (in \$1 increments).

Note: In terms of final payoffs, these games are identical — they both involve Player 1 choosing how to allocate \$10 between Player 1 and Player 2 (in \$1 increments).

But also note: In the second game, there are two ways to choose (\$10,\$0).

Results from Lazeer, Malmendier & Weber (2012)

Social Preferences: Summary

- ▶ Clearly people are not pursuing pure self-interest.
 - ▶ But note: This doesn't mean they don't care about their own outcomes.
 - ▶ If you had to write down a very simple model, pure self interest is probably close enough.
 - ▶ ...that's why this is a 400-level class; the nuance isn't entirely necessary, but you'd otherwise be unable to handle a lot of data and real-world experience.
- ▶ A variety of social motives have been suggested: altruism, inequity aversion, utilitarianism, maximin, reciprocity, concern withdrawal....
- ▶ No single theory can explain all the data — still lots to learn.

Exercise

Suppose that we collect data on how people play the two games depicted on the board (using techniques as in Charness & Rabin). As depicted below, 30% of subjects chose B2 in Game A, while $Y\%$ of subjects chose B2 in Game B. Note: in both games, Player B's payouts are listed **second**.

- (a) If everyone has inequity aversion (with differing degrees of concern), what would we expect for Y ?
- (b) If everyone has social welfare preferences (with differing degrees of concern of social welfare), what would we expect for Y ?
- (c) If everyone has social welfare preferences (with differing degrees of concern of social welfare), and also has concern withdrawal, what would we expect for Y ?
- (d) If everyone has social welfare preferences (with differing degrees of concern of social welfare), and also has strong reciprocity (negative and positive), what would we expect for Y ?