

Introduction to Social Preferences_⇒

- Social preferences can be defined as “ways that people’s utility depend *directly* on the well-being, motives, and beliefs of others” ._⇒
- Among the most famous passages in economics is from Adam Smith’s *The Wealth of Nations*:_⇒
 - 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.

↪

- Research recognizing the centrality of self-interest without ignoring social preferences, illustrated by Dawes and Thaler (1988): \Rightarrow
 - “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. \Rightarrow The box has just a small slit, so money can only be put in, not taken out. \Rightarrow Also, the box is attached to the table, so no one can (easily) make off with the money. \Rightarrow We think that the farmers have just about the right model of human nature. \Rightarrow They feel that enough people will volunteer to pay for the fresh corn to make it worthwhile to put it out there. \Rightarrow The farmers also know that if it were easy enough to take the money, someone would do so.”



Social Preferences

While moving now to the 'field', most evidence gathered from lab. \Rightarrow

- Legitimate concerns about extrapolability from one (campus *or* non-campus) setting to any other (campus or non-campus) setting. \Rightarrow

Note: Ignoring social preferences requires more than (wrong) claim that self interest **always** much more powerful. \Rightarrow

- Rather, *also* requires \Rightarrow
 - No situations where actors can significantly affect well-being of others with only small effect on own well-being. \Rightarrow
- Note: Reference Dependence massively important in social prefs \Rightarrow
 - We ignore.



Social Preferences

I am attitudy and puzzled by much of the research.⇒

- A huge industry in experiments on social preferences.⇒
- Most active area of both experimental and behavioral economics.⇒
- But not the most successful.⇒
 - Little integration of any lab lessons into mainstream economics.⇒
 - Movement towards conceptual tightness disappointingly slow.⇒
 - Movement towards serious empirical conclusions (including acknowledging ones staring us in the face) distressingly slow.⇒
 - And 100% self-interest not as far off the mark as many other assumptions behavioral economists seek to modify, so benefits lower.

↪

Social Preferences

More attitude: \Rightarrow

- Scientifically sinful to not disentangle obvious confounds when they *are* obvious, and when they are easy to disentangle. \Rightarrow
- Multiple theories for the same data rather than creating data to narrow down theories has been egregious. \Rightarrow
 - Even when data needed is really really really *really* easily created. \Rightarrow
- Part of reason: people sense lab data missing something of people's intuition about the world. \Rightarrow
 - Can't get envy in lab...mimic it with other social preferences! \Rightarrow
 - But that suggests not to use the lab—it does not justify adding confounds to the lab so that the data match your intuition.

↪

Two examples, & competing hypotheses. \Rightarrow

- Ultimatum game and Prisoner's Dilemma \Rightarrow

	Accept	Reject
Share	5,5	5,5
Grab	8,2	0,0

	C	D
C	4,4	0,5
D	5,0	1,1

- Variants of these two games dominated much of the earlier social-preferences literature ... \Rightarrow
 - Despite being awful games to do science of social preferences with. \Rightarrow
- Now we have Dictator and Trust games!

\rightarrow

Social Preferences

Consensus empirical fact: people often (sacrifice money to) reject lop-sided offers in ultimatum game. \Rightarrow Why? \Rightarrow

- H1a: failed, mischievous, or lazy selfishness. \Rightarrow
- H2a: punishing obnoxious/unfair behavior. \Rightarrow
- H3a: hate coming out behind random other subject. \Rightarrow

Consensus empirical fact: people often cooperate with others in one-shot prisoner's dilemma, but *only if* they think partner will. \Rightarrow Why? \Rightarrow

- H1b: mistargeted attempt at repeated-game selfishness. \Rightarrow
- H2b: positive reciprocity, rewarding other. \Rightarrow
- H3b: implementing equitable outcome. \Rightarrow

H1a/b: great pedigree, popularity, and wrongness. \Rightarrow

H2a, *H2b*, *H3a*, and *H3b* all reasonable. \Rightarrow

- And could all be true



Social Preferences

Laboratory evidence has inspired many models of “social preferences”. \Leftarrow

- Three classifications, “distributional preferences”, “intentions-based preferences,” and “other belief-based preferences”. \Leftarrow
- Note: economists prone towards particular assumption: altruism, either targeted (children) or need-based (charity). \Leftarrow
 - Both are majorly right. Economists should and do study targeted altruism, and nobody objects to the investigation of charity, etc. \Leftarrow
- But lab research has emphasized fairness, reciprocity, etc.

↪

Distributional Preferences \Rightarrow

- Represent Person 0's preferences by $U_0(\pi_0, \pi_1, \pi_2, \dots)$, \Rightarrow
 - where π_k is Person k 's "material" utility/payoffs. Begin with: \Rightarrow
- "Disinterested distributional preferences": Allocations people choose for others, when their choices do not affect their own outcome?

\rightarrow

Social Preferences

Formally, $W_0(\pi_1, \pi_2, \dots)$. Examples of (extreme) preferences: \Rightarrow

- Surplus-maximizing: $W_0 = \sum_k \pi_k$. \Rightarrow
 - (If π_k really “material hedonic return”, then this is utilitarianism.) \Rightarrow
- Rawlsian/maximin: $W_0 = \text{Min}\{\pi_k\}_k$ \Rightarrow
- Egalitarian: $W_0 = -\sum_k (\pi_k - \bar{\pi})^2$ \Rightarrow

Note: \Rightarrow

- Rawlsian preferences have a form of aversion to inequality, but they are monotonically increasing. \Rightarrow
- “Egalitarian preferences” are qualitatively more extreme dislike of unequal outcomes; you’d actually lower people’s payoffs.) \Rightarrow

These are of course unrealistically extreme forms.

\curvearrowright

Social Preferences

“Non-disinterested” distributional preferences: $\Rightarrow U_0(\pi_0, \pi_1, \pi_2, \dots) =$

$$(1 - k - l)\pi_0$$

$$+ (k \cdot W_0(\pi_0, \pi_1, \pi_2, \dots))$$

$$+ (l \cdot D_0(\pi_0 - \pi_1, \pi_0 - \pi_2, \dots)),$$

where $k, l, k + l \in [0, 1]$. \Rightarrow

- Of course, components are interpretation; don't observe separately. \Rightarrow
- What are they? \Rightarrow
- Answer: \Rightarrow self interest, \Rightarrow disinterested principles \Rightarrow “social comparison”

Social Preferences

CR's simplified parameterization of two-person preferences, based on FS. \Rightarrow

- $U_B(\pi_A, \pi_B) \equiv \rho\pi_A + (1 - \rho)\pi_B$ when $\pi_B \geq \pi_A$. \Rightarrow
- $U_B(\pi_A, \pi_B) \equiv \sigma\pi_A + (1 - \sigma)\pi_B$ when $\pi_B \leq \pi_A$. \Rightarrow

Exaggerates the kinkiness of preferences. \Rightarrow

- This is largely for simplification. \Rightarrow
- Surely generally true that $\rho \geq \sigma$. \Rightarrow
- But the rest up for empirical grabs.

\rightarrow

Social Preferences

A bunch of (not necessarily mutually-exclusive) examples: \Rightarrow

- $\rho = \sigma = 0 \quad \Rightarrow$ pure self interest. \Rightarrow
- $\rho = 1, \sigma = 0 \quad \Rightarrow$ disinterested pure Rawlsian \Rightarrow
- $\rho = \sigma = \frac{1}{2} \quad \Rightarrow$ disinterested surplus maximizer \Rightarrow
- $1 \geq \rho \geq \sigma \geq 0 \quad \Rightarrow$ Social-Welfare Preferences \Rightarrow
- $1 \geq \rho \geq 0 \geq \sigma \quad \Rightarrow$ inequity averse \Rightarrow
- $\rho \geq 1 \geq 0 \geq \sigma \quad \Rightarrow$ egalitarian? \Rightarrow
- $0 \geq \rho \geq \sigma \quad \Rightarrow$ competitive? \Rightarrow

\rightarrow

Lab Evidence on Social Preferences \Rightarrow

- In all examples, proportions choosing each of two money combinations (usually U.S. pennies or Spanish pesetas) are drawn underneath the amounts. \Rightarrow
 - From Kritikos & Bolle (2001), Charness & Grosskopf (2001), and Charness & Rabin (2002,2005): \Rightarrow
- Caution: \Rightarrow
 - These slides somewhat dated \Rightarrow
 - More evidence on some of this \Rightarrow
 - (Although remarkably little)



Some Lab Data on Distributional Prefs \Rightarrow

Evidence on disinterested distributional preferences? \Rightarrow

- There is virtually no \$-stakes evidence on disinterested preferences! \Rightarrow
- From Charness and Rabin (2002): \Rightarrow

C chooses (A,B) allocation of $(400,400)$ vs. $(750,375)$
.46 vs. .54 \Rightarrow

C chooses (A,B) allocation of $(400,400)$ vs. $(1200,0)$
.82 vs. .18 \Rightarrow

Keep these in mind. \Rightarrow

- Especially the first one \Rightarrow
- Disinterested care a lot about "equity"

Evidence on Non-Disinterested Distributional Preferences?

- Quick, selective examples meant to illustrate. \Rightarrow
- (But meant not to be misleading.) \Rightarrow
- To ask purely distributional preferences, free of reciprocity, consider first only “dictator” games. \Rightarrow
 - My attitude begins here ... \Rightarrow
 - disentangle rather than confound.

↗

Social Preferences

What is the evidence on ρ ? \Rightarrow

B chooses (A,B) allocation of	(200,700)	vs.	(600,600)	\Rightarrow
Old C-R:	.27		.73	

B chooses (A,B) allocation of	(0,800)	vs.	(400,400)	\Rightarrow
Old C-R	.78		.22	
New C-R	.56		.44	
New C-R with requests by A	.45		.55	

My impressions from the accumulated experimental evidence: \Rightarrow

- The average or median ρ is about .4. \Rightarrow
- About 10% of subjects have $\rho < 0$.

\rightarrow

Evidence on σ ? \Rightarrow

- Shockingly (unforgivably) little evidence on this disentangled from negative reciprocity. \Rightarrow
 - Do people Pareto-damage when other has done nothing wrong? \Rightarrow
- Warning: might be misleading, because compared to other experiments, relatively little “Pareto damage”. \Rightarrow

\Rightarrow

Social Preferences

Small sacrifice to avoid coming out behind? \Rightarrow

B chooses (A,B) allocation of (1200,625) vs. (600,600)
C-G 88% 12% \Rightarrow

B chooses (A,B) allocation of (750,400) vs. (375,375)
C-R 77% 23% \Rightarrow

That 23% is the largest % I know of for reciprocity-free strict Pareto-damaging sacrifice, across all experiments ever run.

- Recent evidence, still absorbing, finds more. \Rightarrow
- (And certainly also more evidence of less.)

\curvearrowright

Social Preferences

Significant sacrifice to avoid coming out behind? \Rightarrow

B chooses (A,B) allocation of (4,1) vs. (0,0)
K-B 88% 12% \Rightarrow

B chooses (A,B) allocation of (3,2) vs. (0,0)
K-B 100% 0% \Rightarrow

B chooses (A,B) allocation of (800,200) vs. (0,0)
C-R 100% 0%

\curvearrowright

Social Preferences

Small sacrifice to come out *further* behind? \Rightarrow

B chooses (A,B) allocation of $(625,625)$ vs. $(1200,600)$
 $C-G$ 33% 67% \Rightarrow

B chooses (A,B) allocation of $(400,400)$ vs. $(750,375)$
 $C-R$ 55% 45% \Rightarrow

- Recall 1: Only 23% choose $(375,375)$ over $(750,400)$. \Rightarrow
- Recall 2: 46% of *disinterested* choose $(400,400)$ over $(750,375)$. \Rightarrow
- **So:** more people sacrifice to come out **behind** when it really helps other than will sacrifice to avoid coming out behind! \Rightarrow
- **And:** most who do not sacrifice to help: \Rightarrow are refusing out of sincere Rawlsian motives, not because selfish or because hate coming out behind.



Social Preferences

Costlessly take \$ from other to avoid behind? \Rightarrow

<i>B chooses (A,B) allocation of</i>	$(X,0)$	vs.	$(0,0)$	
<i>K-B</i>	75%		25%	\Rightarrow
<i>B chooses (A,B) allocation of</i>	$(900,600)$	vs.	$(600,600)$	
<i>C-G</i>	67%		33%	\Rightarrow
<i>B chooses (A,B) allocation of</i>	$(750,400)$	vs.	$(400,400)$	
<i>C-R</i>	68%		32%	\Rightarrow
<i>B chooses (A,B) allocation of</i>	$(2000,400)$	vs.	$(400,400)$	
<i>C-R</i>	82%		18%	\Rightarrow

- More than I ever expected... $\frac{1}{3}$ taking \$ away from other. \Rightarrow
 - But recall: behindness aversion says motive super-strong here.



Social Preferences

Crude summary accumulated experimental evidence: \Rightarrow

- About 30% $\sigma < 0$, about 70% $\sigma > 0$. \Rightarrow
- Median $\bar{\sigma} > 0$, but very few $|\sigma| \gg 0$. \Rightarrow
- Including very few $\sigma \ll 0$. \Rightarrow

Update on original examples and hypotheses? So far: \Rightarrow

- Little indication that rejections in the ultimatum game have much to do with “behindness aversion”. \Rightarrow
 - Maybe it is all about negative reciprocity instead. \Rightarrow
- But lots of evidence of sharing when ahead ... very consistent with PD being from “inequity aversion”. \Rightarrow
 - Does not preclude positive reciprocity.

\curvearrowright

Intentions-Based Preferences \Rightarrow

- People may care not *just* about outcomes, but with rewarding and punishing good and bad behavior. \Rightarrow
 - In bilateral context, we think of this as reciprocity. \Rightarrow
- I emphasize “**just**” to make clear: \Rightarrow
 - conceptually incoherent to have preferences that are just about reciprocity ... \Rightarrow
 - people **must** have some notion of good and bad outcomes in order to be reciprocal about anything. \Rightarrow
 - Rabin (1993), Dufwenberg-Kirchsteiger embed bad distributional \Rightarrow
 - Falk-Fischbacher, Charness-Rabin combine more serious distributional preferences with reciprocity. \Rightarrow

↪

How to test for the role of intentions? \Rightarrow

- Very simple method: \Rightarrow
 - With common knowledge to the players, Player A chooses between outcome X and giving Player B the choice from $\{Y, Z\}$. \Rightarrow
 - Player B's preferences between Y and Z across situations, or when Player A has no choice, reflects B's distributional preferences. \Rightarrow
 - But the way Player B's preferences between Y and Z depend on changes in X reflect his reciprocal/intentions-based preferences.

\rightarrow

Social Preferences

Begin with the dark side: What induces “Pareto-damaging” behavior? \Rightarrow

- Behavior that hurts some or all without helping anybody. \Rightarrow

<i>A chooses</i>	<i>or</i>	<i>lets B choose</i>	$(800,200)$	<i>vs.</i>	$(0,0)$	
<i>No choice</i>			100%		0%	\Rightarrow
<i>fairer than</i>		$(800,200)$	81-92%		8-19%	

- Note: discussing only B behavior. \Rightarrow
 - A behavior in papers. \Rightarrow
- Reminder: these data not typical evidence in the literature, where more Pareto damage is typically observed.

\rightarrow

Social Preferences

<i>A chooses</i>	<i>or let</i>	<i>B choose</i>	$(750,400)$	vs.	$(375,375)$	
<i>No choice</i>			77%		23%	
<i>fairer than</i>	$(750,400)$		71%		29%	\Rightarrow
$(400,750)$			80%		20%	

Increase is statistically significant, but not large. Note that seems that B does not punish A for not wanting to be on the short end herself.

↪

Social Preferences

When B makes choices where \$ for A at stake, but not \$ for B, how respond to goodness, badness? \Rightarrow

<i>A chooses</i>	<i>or</i>	<i>lets B choose</i>	$(750,400)$	vs.	$(400,400)$	
<i>No choice</i>			$\approx 60\%$		$\approx 40\%$	
$(550,550)$			$\approx 55\%$		$\approx 45\%$	\Rightarrow
$(750,0)$			94%		6%	

- What is going on? \Rightarrow
 - 1st vs. 2nd row? \Rightarrow
 - 1st vs. 3rd? \Rightarrow
- Surprisingly little punishment, even when free. \Rightarrow
 - A puzzle. \Rightarrow
- **But behindness aversion virtually vanishes when A has been kind.**



Crude summary of accumulated experimental evidence: \Rightarrow

- Not much Pareto-damage without reciprocity \Rightarrow
- Increase in Pareto-damage by B if A is mean/unfair \Rightarrow
- Stronger indication of diminishing Pareto-damage if A behaves nicely.

\Rightarrow

Social Preferences

How does good and bad behavior by one player affect the other player's inclination to engage in helpful sacrifice? \Rightarrow

	<i>A chooses</i>	<i>or lets B choose</i>	$(750, 375)$	vs.	$(400, 400)$	
	<i>No choice</i>		46%		54%	
\Rightarrow	$(750 \pm 50, 0)$		37%		63%	\Rightarrow
	$(550, 550)$		11%		89%	

- Comment on 1st vs. 2nd line: \Rightarrow
 - Whoa! \Rightarrow
- Lots and lots of data since we didn't believe the results. \Rightarrow
 - It's robust. \Rightarrow
 - Statistically significantly in opposite direction as positive reciprocity.

\rightarrow

Social Preferences

Accumulated evidence: \Rightarrow

- Lots of helpful sacrifice.
- *Not* increased by other's good behavior. \Rightarrow
- But withdrawn if other misbehaves. \Rightarrow

Some experiments find positive reciprocity, but very few. \Rightarrow

- My bet: meta-analysis would show people that positive reciprocity in sense of behaving better towards somebody if she has been good is approximately **zero** in the lab. \Rightarrow
- The “concern withdrawal”, sort of in between positive reciprocity and negative reciprocity as conventionally conceived, seems robust.

↪

Social Preferences

- Explanations, by Bozos like Rabin (1993), that cooperation in laboratory PD is from positive reciprocity, is wrong. \Rightarrow
- Similarly: massive experiments on “trust” seem really to be about ρ — little indication that trust is rewarded independent of desire to share. \Rightarrow
- I believe in positive reciprocity \Rightarrow
 - and have thoughts why missing from lab. \Rightarrow
 - but it is not in evidence whatsoever in the lab.

↪

Social Preferences

- Approximately: \Rightarrow
 - In UG: All negative reciprocity, no behindness aversion. \Rightarrow
 - In PD: All equity, no positive reciprocity. \Rightarrow
- Over-strong statements? \Rightarrow
 - Stark statements to counter-balance the apparent triumph of priors and theory and pet models over easy-to-see and overwhelming empirical evidence. \Rightarrow
- The point is **not** a horse race, \Rightarrow
 - not a claim about which motives bigger. \Rightarrow
- Estimates in lab on behindness aversion and positive reciprocity. \Rightarrow
- Both are (slightly) backwards.

↪

Back to belief-based utility? \Rightarrow

- Reciprocity models (and cousins/descendants a la Levine "Spite") are clearly belief-based. \Rightarrow
- But so are the new models based on image. \Rightarrow
- Clearly right direction \Rightarrow
 - Even if some of them formulated in complacent terms. \Rightarrow
 - And must be "seeded" by distributional assumptions. \Rightarrow
 - Because all models must.

\rightarrow

Social Preferences

Great early experiment by Dana, Weber, and Kuang (2007): \Rightarrow

- (other, self) of $A \equiv (\$5, \$5)$ vs. $B \equiv (\$1, \$6)$. \Rightarrow
 - 26% choose $B \equiv (\$1, \$6)$. \Rightarrow
- (other, self) of $C \equiv (\$1, \$5)$ vs. $D \equiv (\$5, \$6)$. \Rightarrow
 - Presumptively, 100% choose D . \Rightarrow
 - And, in hypothetical treatment, 100% in fact chose it. \Rightarrow
- Then: people told to choose either \$5 or \$6 for self. \Rightarrow
 - But 50% chance $(\$ _, \$5)$ vs. $(\$ _, \$6)$ is A vs. B . \Rightarrow
 - 50% chance it is C vs. B .

\rightarrow

Social Preferences

- 37% choose to Not reveal, take \$6.⇒
- 7% choose to Not reveal, take \$5.⇒
- 56% choose to Reveal.⇒
- If revealed and saw $(\$5, \$5)$ vs. $(\$1, \$6)$,⇒
 - 25% choose $(\$1, \$6)$ ⇒
- If revealed and saw $(\$1, \$5)$ vs. $(\$5, \$6)$,⇒
 - 90% choose $(\$5, \$6)$. (So 10% are either silly or nasty)⇒
- And so if choice turned out to be $(\$5, \$5)$ vs. $(\$1, \$6)$, 63% (rather than 26%) ended up choosing $(\$1, \$6)$, either by choosing it after seeing it, or by choosing not to see it.

↪

Social Preferences

Interpretation? \Rightarrow

- “Moral Wiggle Room”. Authors (incessantly) point out cannot explain by distributional preferences. Why not reveal payoffs? \Rightarrow
- One interpretation: self-image preservation. \Rightarrow
- Grossman (2008) and Lazear, Malmendier, and Weber (2007), Tadelis, etc.: variants of people seeming to have some variant of belief-based social preferences.
 - Shame? \Rightarrow
 - Social signaling? \Rightarrow
 - Self signaling? \Rightarrow
- People care how obvious it is to themselves and to others that they are being selfish or virtuous.



Final point \Rightarrow

- “Wiggle room \rightarrow non-selfish behavior may be less important than it appears” is not an awful intuition. \Rightarrow
- But it’s not as compelling as may seem. \Rightarrow
- A world where people convince themselves they are right to justify selfish behavior may be *more* different from simple selfishness than is simple simple social preferences. \Rightarrow
- Bargaining selfish people vs. bargaining altruists vs. bargaining self-righteous people?

\rightarrow

Modeling Reciprocity? Hard. \Rightarrow

- I mean hard in both the sense of technically difficult and in terms of really finding robust models. Other- or self-signaling use signaling models (and moral wiggle room and related experiments seem to require weirdness; Rabin (1995) “moral rules” vs. “moral preferences” with weird assumption about treatment of information. Reciprocity uses “psychological games” by Geanakoplos, Pearce, and Stacchetti (1989), or related approaches, such as Levine’s signaling approach or steady-state equilibrium a la Charness & Rabin. \Rightarrow
- Technical issues in modeling reciprocity: Both as psychological reality and as modeling strategy, not just outcomes, nor even cleverly defined to include game, but beliefs. \Rightarrow
- How might we capture the role of volition and intentions in models of social preferences?



Social Preferences

Claim: Often no $U_j(\pi_i, \pi_j)$ can capture preferences. \Leftarrow

- Nor, to slightly generalize, can capture by writing payoffs solely function of outcomes. \Leftarrow
- Point is not that you can't do in terms of own payoff. Duh. The point is that it is not a function of everybody's payoff ... but different things. In any framework, psychological games of type-signaling games, need beliefs in there. \Leftarrow
- Think about the following game. Fully hypothetical; no data. Going to 3-person for ease of illustration.

↪

Social Preferences

- Want you to think about the following game (using intuition, like I am, since nobody has run the game to my knowledge). Going to 3-person for ease of illustration. \Rightarrow

	1 goes L		1 goes R		
	2 chooses		3 chooses		
	L	R	L	R	
1's Payoff	10	12	11	5	\Rightarrow
2's Payoff	10	12	5	5	
3's Payoff	10	0	5	5	

What do you think Player 3 will do if Player 1 goes R? What are the issues? \Rightarrow

- If Player 3 thinks Player 2 would have gone L? \Rightarrow
- If Player 3 thinks Player 2 would have gone R? \Rightarrow
- What is not right about those questions?