

# Eliciting Moral Preferences: Theory and Experiment

Roland Bénabou, Armin Falk, Luca Henkel and Jean Tirole

Paper available via web page

ESA World Meeting – September 2020

# Introduction

- What can be learned about a person's or a population's moral preferences from observing their choices, including in experiments?
  - ▶ How should this be used to inform policy, or to maximize prosocial actions?
  - ▶ How to interpret behaviors that seem deontologically rather than consequentially motivated: refusing tradeoffs that involve harm to others, assigning infinite price to “sacred values” such as life, freedom, dignity?
- Show how, whenever image concerns are present, the answers depend crucially on **how choices are elicited / contributions solicited**:
  - ▶ **Single** (or, separate) decisions, vs. **multiple** simultaneous decisions; e.g., yes/no to an offer, versus stating a willingness to pay
  - ▶ Ex-ante commitments under **uncertainty**, vs. known, **ex-post** choices; e.g., random realized situation, random implementation

# Concrete Setting and Application

- Use model & experiment to study and compare properties of two most commonly used revealed-preference methods:
    - ▶ Direct elicitation (*DE*)
    - ▶ Multiple-price list (*MPL*)
  - Compared to *DE*, *MPL* features **multiple** decisions, of which only one is implemented for real, **at random**
  - In standard situations (e.g., for non-moral decisions), we know both schemes give the **same, and correct**, answer. For instance:
    - ▶ Ask people in a population to make a *DE* choice, each one at different price
    - ▶ Ask each person the same *MPL* choice question
- ⇒ Get same distribution of outcomes, estimate same distribution of preferences

# Key Results

- 1 As soon as image concerns are present, *DE* and *MPL* give different answers
  - ▶ Unrelated to risk aversion: tradeoffs now differ in expected value
- 2 Gap between results *varies* with the importance of image concerns (interaction), not just in *magnitude*, but even in *sign*! At any given price:
  - ▶ *DE* will generate more prosocial behavior than *MPL* when image concerns are weak (but positive)
  - ▶ *MPL* will generate more prosocial behavior than *DE* when image concerns are strong
- 3 Image-minded consequentialists will display *Kantian-like* price insensitivity much more readily under MPL than under DE
- 4 Results due interplay of *three general effects*, also at work in public-goods contributions mechanisms sharing key features with DE / MPL.
  - ▶ *Discouragement effect, cheap-talk effect, and cheap-act effect*
- 5 Model's most distinctive prediction: “*crossing pattern*” between *DE* and *MPL* contributions, as image concern go from weak to strong
- 6 Test it in an experiment: taking money or “*Saving a Life*”

# Model

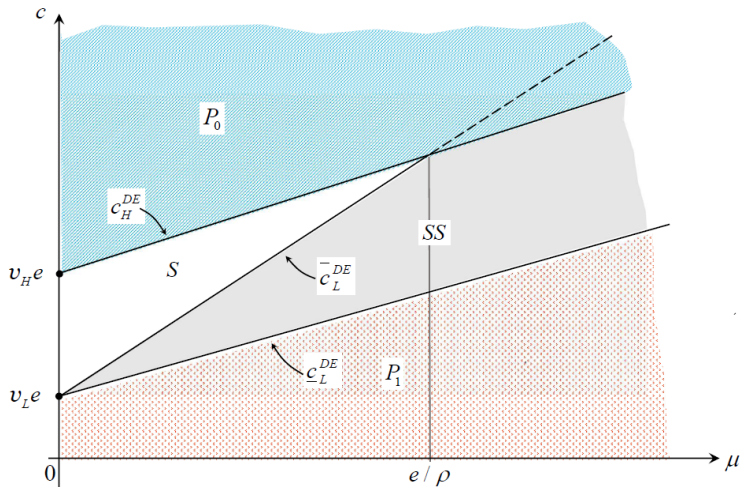
- Choice: engage in moral behavior ( $a = 1$ ) or act selfishly ( $a = 0$ )
  - ▶  $a = 1$  involves personal cost  $c$  but generates positive externality  $e$
- Agents differ in their motivation to act morally:
  - ▶ High type  $v_H e$ , with prob  $\rho$ , Low type  $v_L e$ , with prob  $1 - \rho$ ;  $v_H > v_L \geq 0$
- Final utility for type  $\tau = L, H$ :

$$U_\tau(a) = (v_\tau e - c)a + \mu E[v|a, \text{choice conditions}]$$

- $\mu \geq 0$ : strength of self or/and social image concerns. Image / esteem based on agent's expected type, conditional on action  $a$  and choice conditions
- Situation, even experiment, is now a signaling game  $\Rightarrow$  behavior reflects not just individual preferences, but equilibrium
  - ▶ Pareto dominance as selection criteria in case of multiple equilibria

# Behavior under Direct Elicitation

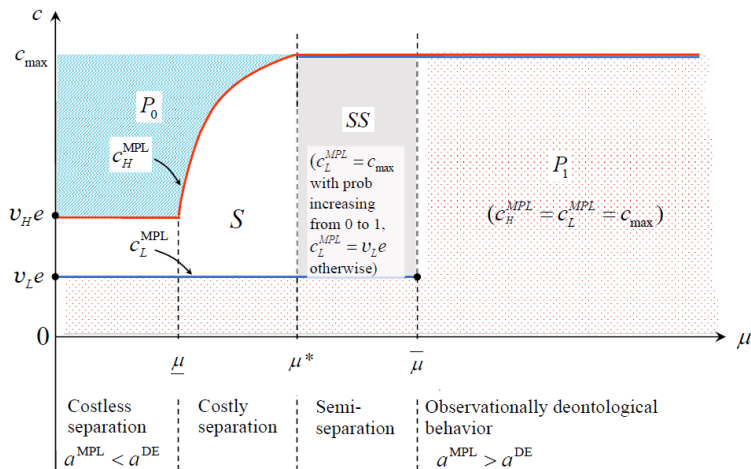
- Agents face choice  $a \in \{0, 1\}$ , for given value  $c \in [0, c_{\max}]$



- $P_0$ : pooling at  $a_H = a_L = 0$ ;  $S$ : separation,  $a_H = 1, a_L = 0$ ;  $SS$ : semi-separation:  $a_H = 1, a_L \in (0, 1)$ ;  $P_1$ : pooling at  $a_H = a_L = 1$ .

# Behavior under Multiple-Price List

- Agents state maximum level of  $c$  to take  $a = 1$  (WTP)
- Actual  $\tilde{c}$  drawn from  $G(\tilde{c})$  on  $[0, c_{\max})$ , implement  $a = 1$  at cost  $\tilde{c}$  iff  $\tilde{c} \leq c$



- $P_0$ : pooling at  $a_H = a_L = 0$ ;  $S$ : separation,  $a_H = 1, a_L = 0$ ;  $SS$ : semi-separation:  $a_H = 1, a_L \in (0, 1)$ ;  $P_1$ : pooling at  $a_H = a_L = 1$ .

# Intuition: Three Key Effects from DE to MPL I

- ① **Discouragement effect:** because it reveals multiple decisions at the same time, *MPL* raises the cost to the Low type of mimicking the High type:
  - ▶ Say,  $v_{Le} = 50$ ,  $v_{He} = 75$ . Low type might be willing to pool at *DE* price of  $c = 60$ , but under *MPL* would have to be willing to pool up to 75.
  - ▶ If  $\mu$  is positive but not very large, not worth it  $\Rightarrow$  will simply state WTP of 50, and thus does not contribute when  $\tilde{c} = 60$  is drawn
  - ▶ This effect dominates at low  $\mu > 0 \Rightarrow$  *DE* induces more prosocial decisions than *MPL*
- ② **Cheap-talk effect:** under *DE*, if say yes to  $c$ , then  $(e, -c)$  occurs for sure. Under *MPL*, if state WTP of  $c$  there is a probability  $1 - G(c)$  that won't be "called on it", and neither  $e$  nor  $-c$  occurs.
  - ▶ This effect tends to induces more prosocial decisions (especially by the low type) under *MPL*, relative to *DE*
  - ▶ However, as  $\mu$  rises it weakens and ultimately vanishes, as the cutoff  $c_{\tau}^{MPL}$  rises toward  $c_{\max}$ , driving the probability of implementation toward 1.



# Intuition: Three Key Effects from DE to MPL II

- **Cheap-act effect:** under DE, if say yes to  $c$ , then I pay  $c$  for  $e$ . Under *MPL*, if state WTP of  $c$  and am “called on it”, will pay some random  $\tilde{c} \leq c$ .

Because  $c - E_G[\tilde{c} | c \leq \tilde{c}] > 0$ , this effect also tends to induces more prosocial decisions under *MPL*, relative to *DE*.

- ▶ Moreover, for the experimentally standard uniform distribution, and more generally for any distribution  $G(\tilde{c})$  satisfying *MLRP*, the previous difference increases with  $c$ .
- ▶ Therefore, as  $\mu$  rises, pushing up all cutoffs, this effect strengthens. It is thus the one that dominates at high  $\mu$ .
- Intermediate  $\mu$ 's : all three effects operate, not much can be said in general
  - ▶ Paper derives a sufficient condition for single crossing of aggregate contributions under *DE* vs. *MPL* in the case of uniform  $G$ .

# Main Result: Comparing DE and MPL

## Proposition (interactions and reversal)

For each type (hence also on average):

- ① For any  $c \in [0, c_{\max}]$ ,  $a_{\tau}^{\text{MPL}}(c, \mu)$  and  $a_{\tau}^{\text{DE}}(c, \mu)$  *coincide at  $\mu = 0$* , then both *increase (weakly)* as  $\mu$  rises, reaching 1 for  $\mu$  large enough.

- ② For all  $\mu \in (0, \underline{\mu})$ ,

$$a_{\tau}^{\text{DE}}(c, \mu) \geq a_{\tau}^{\text{MPL}}(c, \mu),$$

with strict inequality for  $c \in (v_L e, \underline{c}_L^{\text{DE}}(\mu))$  and  $c \in (v_H e, \underline{c}_H^{\text{DE}}(\mu))$ , both nonempty.

- ③ For all  $\mu \geq \bar{\mu}$ ,

$$a_{\tau}^{\text{DE}}(c, \mu) \leq a_{\tau}^{\text{MPL}}(c, \mu),$$

with strict inequality for  $c \in (\underline{c}_L^{\text{DE}}(\mu), c_{\max})$ , which is nonempty whenever  $\mu \in (\bar{\mu}, \mu^{**})$ .

# Empirical Tests

- **Hypothesis 1:** For both *DE* and *MPL*, total contributions increase in  $\mu$
- **Hypothesis 2:** For low  $\mu_L > 0$ , total contributions are higher under *DE* than under *MPL*
- **Hypothesis 3:** For high  $\mu_H$ , total contributions are higher under *MPL* than under *DE*
- **Corrollaries:**
  - ▶ **Differential image sensitivity:** as  $\mu$  changes from  $\mu_L$  to  $\mu_H$ , contributions rise by more under *MPL* than under *DE*
  - ▶ **Observationally deontological behavior:** at  $\mu_H$ , more people will choose the moral action “whatever it costs” , i.e. up to the highest price  $c_{\max}$  under *MPL*, than under *DE*. Different estimated fractions of “Kantians”.

# Experiment: Saving a Life

- **Choices:**

- ▶ Moral action ( $a = 1$ ) : induce a 350€ donation that, in expectation / on average, will save one patient from death by tuberculosis. Major  $e \gg 0$
- ▶ Selfish action ( $a = 0$ ) : take money for oneself.  
Amount  $c$ , can range from 10 to 200€

- **High stakes:**

- ▶ Subjects provided with detailed, verifiable (on site) evidence of death risk for tuberculosis patients in India, effectiveness of treatment, track record of NGO doing it (Operation ASHA), expected value calculation

- **Treatments:**  $2 \times 2$  between-subjects design, varying both:

- ▶ **Elicitation method:** Direct elicitation (*DE*) vs. Multiple-price list (*MPL*)
- ▶ **Level of image concerns,  $\mu$**  : choices kept private (*Low Image*), or made publicly visible & morally salient (*High Image*)

# Decision Screens

## Your Decision

Please click here to be reminded of the precise meaning of 'saving a life'

| Option A            |                       |                       | Option B                             |
|---------------------|-----------------------|-----------------------|--------------------------------------|
|                     | A                     | B                     |                                      |
| I save a human life | <input type="radio"/> | <input type="radio"/> | I choose 100 € as payment for myself |

Confirm decision

## Your Decisions

Please click here to be reminded of the precise meaning of 'saving a life'

| Option A            |                       |    |                       | Option B                             |
|---------------------|-----------------------|----|-----------------------|--------------------------------------|
|                     | A                     |    | B                     |                                      |
| I save a human life | <input type="radio"/> | 1  | <input type="radio"/> | I choose 0 € as payment for myself   |
| I save a human life | <input type="radio"/> | 2  | <input type="radio"/> | I choose 10 € as payment for myself  |
| I save a human life | <input type="radio"/> | 3  | <input type="radio"/> | I choose 20 € as payment for myself  |
| I save a human life | <input type="radio"/> | 4  | <input type="radio"/> | I choose 30 € as payment for myself  |
| I save a human life | <input type="radio"/> | 5  | <input type="radio"/> | I choose 40 € as payment for myself  |
| I save a human life | <input type="radio"/> | 6  | <input type="radio"/> | I choose 50 € as payment for myself  |
| I save a human life | <input type="radio"/> | 7  | <input type="radio"/> | I choose 60 € as payment for myself  |
| I save a human life | <input type="radio"/> | 8  | <input type="radio"/> | I choose 70 € as payment for myself  |
| I save a human life | <input type="radio"/> | 9  | <input type="radio"/> | I choose 80 € as payment for myself  |
| I save a human life | <input type="radio"/> | 10 | <input type="radio"/> | I choose 90 € as payment for myself  |
| I save a human life | <input type="radio"/> | 11 | <input type="radio"/> | I choose 100 € as payment for myself |
| I save a human life | <input type="radio"/> | 12 | <input type="radio"/> | I choose 110 € as payment for myself |
| I save a human life | <input type="radio"/> | 13 | <input type="radio"/> | I choose 120 € as payment for myself |
| I save a human life | <input type="radio"/> | 14 | <input type="radio"/> | I choose 130 € as payment for myself |
| I save a human life | <input type="radio"/> | 15 | <input type="radio"/> | I choose 140 € as payment for myself |
| I save a human life | <input type="radio"/> | 16 | <input type="radio"/> | I choose 150 € as payment for myself |
| I save a human life | <input type="radio"/> | 17 | <input type="radio"/> | I choose 160 € as payment for myself |
| I save a human life | <input type="radio"/> | 18 | <input type="radio"/> | I choose 170 € as payment for myself |
| I save a human life | <input type="radio"/> | 19 | <input type="radio"/> | I choose 180 € as payment for myself |
| I save a human life | <input type="radio"/> | 20 | <input type="radio"/> | I choose 190 € as payment for myself |
| I save a human life | <input type="radio"/> | 21 | <input type="radio"/> | I choose 200 € as payment for myself |

Confirm decisions

# Manipulating Moral-Image Concerns

- To ensure some minimal social and self-image  $\mu > 0$ , subjects are anonymously paired, will learn partner's choices (benchmarking). Then:

## ► Low Image ( $\mu_L$ ) :

- Experiment is double blind. Use procedure of Barmettler, Fehr, and Zehnder (2012): one subject carries out final payment, without participating in experiment. Self-image still presumably operating.

## ► High Image ( $\mu_H$ ) :

- Subject's choices are publicly **observed and compared** to those of their partners by a committee, upon receiving payment
  - Own and partner's choices **projected on a wall** with subject present, must **read them aloud**.
  - Committee of three, sitting in the room, **evaluates morality** of each choice
    - ★ Morality scores not disclosed, but really given, and subjects know that

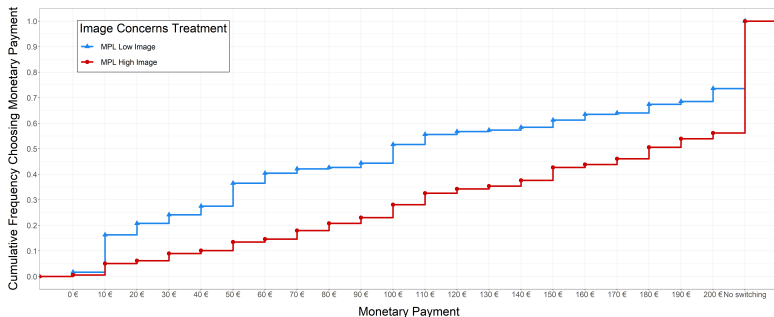
# Procedure

- Bonn Lab: 697 subjects, mostly students, 58% female, mean age = 24.01
- 12€ show-up fee. Receive extensive background information on donation, decisions must pass comprehension test. For each session ( $\approx 20$  subjects) the decisions of one pair are implemented for real.
- Implement, under High and Low Image:
  - ▶ *DE* at (preset) price of  $c = 100\text{€}$
  - ▶ *MPL* with  $\tilde{c}$  uniform over  $0, 10, \dots, c_{\max} = 200\text{€}$ , in increments of  $10\text{€}$
- When comparing the two schemes, do so at same price level:
  - ▶  $\bar{a}^{DE}(100, \mu)$  : fraction who save a life rather than take  $c = 100\text{€}$ , under *DE* (would then have done so under *DE* at any  $c' \leq 100$ )
  - ▶  $\bar{a}^{MPL}(100, \mu)$  : fraction who state  $\text{WTP} \geq 100\text{€}$  under *MPL*, and thus commit to saving a life at any  $\tilde{c} \leq 100$  that may be drawn



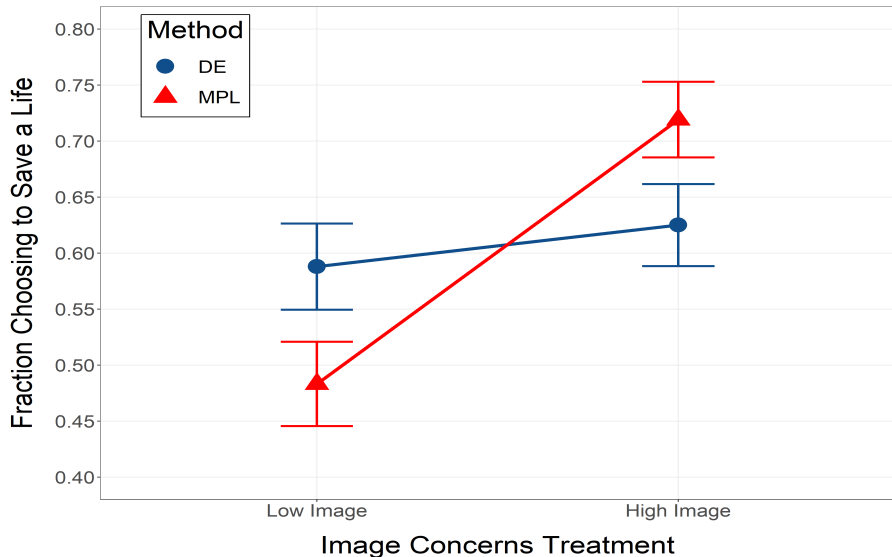
# Hypothesis 1: direct effect of image

- **DE:** 58.8% of subjects choose to save a life (vs. 100€) in *Low Image*, and 62.5% in *High Image*. But difference not significant
- **MPL:**



- CDF from Low Image lies above that from MPL-High Image, for all monetary payments ( $p < 0.001$ , Kolmogorov–Smirnov test). Difference  $> 15\%$  for almost all payments; largest at 60€, of 26%.
- Obs. deontological: 26.4% under  $\mu_L$ , nearly doubles to 48.4% under  $\mu_H$ !

## Hypotheses 2 and 3: interaction and reversal



## Hypotheses 2 + 3: differential image sensitivity

| Dependent variable:    | Choice to Save a Life (vs. 100€) |                     |                            |                     |
|------------------------|----------------------------------|---------------------|----------------------------|---------------------|
|                        | <i>Low Image Concerns</i>        |                     | <i>High Image Concerns</i> |                     |
|                        | (1)                              | (2)                 | (3)                        | (4)                 |
| <i>MPL</i>             | -0.105*<br>(0.054)               | -0.103*<br>(0.053)  | 0.094*<br>(0.050)          | 0.091*<br>(0.050)   |
| Constant ( <i>DE</i> ) | 0.588***<br>(0.038)              | 0.626***<br>(0.049) | 0.625***<br>(0.037)        | 0.622***<br>(0.046) |
| Controls               |                                  | X                   |                            | X                   |
| Observations           | 343                              | 343                 | 354                        | 354                 |
| R <sup>2</sup>         | 0.011                            | 0.077               | 0.010                      | 0.062               |

Robust standard errors in parentheses. Controls include age, gender, income, religiousness, educational level, and high school grade. Significance levels: \* $p < 0.1$ , \*\* $p < 0.05$  and \*\*\* $p < 0.01$ .

# Heterogeneity among subjects

- Use independent measure of altruism (validated in Falk et al. 2018):
  - ▶ “How willing are you to give to good causes without expecting anything in return?”
  - ▶ “Today you unexpectedly received 1,000€. How much of the money would you donate to a good cause?”

⇒ Median split

- Measure is correlated with “saving a life” decision, but independent of treatment

# Heterogeneity among subjects

| Dependent variable: | Choice to Save a Life (vs. 100€) |                     |                     |                     |                       |                     |                     |                     |
|---------------------|----------------------------------|---------------------|---------------------|---------------------|-----------------------|---------------------|---------------------|---------------------|
|                     | Below-median Altruism            |                     |                     |                     | Above-median Altruism |                     |                     |                     |
|                     | Low Image                        |                     | High Image          |                     | Low Image             |                     | High Image          |                     |
|                     | (1)                              | (2)                 | (3)                 | (4)                 | (5)                   | (6)                 | (7)                 | (8)                 |
| MPL                 | -0.187**<br>(0.075)              | -0.187**<br>(0.078) | 0.040<br>(0.075)    | 0.030<br>(0.079)    | -0.032<br>(0.073)     | -0.007<br>(0.072)   | 0.118*<br>(0.068)   | 0.138**<br>(0.066)  |
| Constant (DE)       | 0.512***<br>(0.056)              | 0.611***<br>(0.084) | 0.592***<br>(0.050) | 0.586***<br>(0.072) | 0.663***<br>(0.052)   | 0.591***<br>(0.070) | 0.667***<br>(0.054) | 0.647***<br>(0.062) |
| Controls            |                                  | X                   |                     | X                   |                       | X                   |                     | X                   |
| Observations        | 342                              | 342                 | 342                 | 342                 | 355                   | 355                 | 355                 | 355                 |
| R <sup>2</sup>      | 0.036                            | 0.133               | 0.002               | 0.035               | 0.001                 | 0.101               | 0.017               | 0.109               |

Robust standard errors in parentheses. Controls include age, gender, income, religiousness, educational level, and high school grade. Significance levels: \* $p < 0.1$ , \*\* $p < 0.05$  and \*\*\* $p < 0.01$ .

# Conclusion

- Image concerns **interact** differently with different **elicitation methods**, solicitation schemes. The introduction of **multiple decisions** and **random implementation** give rise to three key effects:
  - ▶ *Discouragement effect*: multiple decisions (WTP) decrease contributions. Dominates at low  $\mu$
  - ▶ *Cheap-act effect*: random cost increases contributions. Dominates at high  $\mu$
  - ▶ *Cheap-talk effect*: operates in middle range.
- **Experimental evidence**: *DE* and *MPL* “crossing” in high-stakes experiment
- **Implications**:
  - ① Caveat for **measurement of moral preferences** (and other reputation-bearing behaviors), whether one is interested in descriptive / predictive questions (how people behave, including from reputation-seeking), or normative ones (how much they truly value public goods, behaviors)
  - ② Caveat about / upper bound on / estimating the proportion of “**Kantians**”
  - ③ Possible applications / extensions to other types of preference elicitation schemes, charitable-contributions solicitations, etc.