

# Ethics of randomized field experiments: Evidence from a randomized survey experiment

Hide-Fumi Yokoo<sup>†</sup>

Hitotsubashi University and Research Institute of Economy, Trade and Industry (RIETI)

April 2023

## Abstract

Using online surveys in Japan, I empirically study the ethical concerns regarding randomized field experiments. Among the six existing experiments, an early childhood intervention is recognized as the most acceptable, while a charitable fund-raising experiment using lotteries is recognized as the least acceptable from an ethical perspective. I find a nonsignificant impact of changing the research methodology from a randomized experiment to an uncontrolled before–after study. However, ethical concerns significantly increase when informed consent is not enough or when subjects are randomly sampled. These findings support an experiment with agreed-upon participants, although it limits the external validity.

*Keywords:* Ethical issues, Field experiments, Online surveys, Randomized controlled trials

JEL classification: C93, D63, O22

---

\* I thank Keitaro Aoyagi, Yuki Higuchi, Yohei Kobayashi, Koichi Kuriyama, Aya Suzuki, Kenji Takeuchi, Fumio Ohtake, Junichi Yamasaki, and participants at a seminar at the Research Institute of Economy, Trade and Industry (RIETI), Hitotsubashi, as well as conference participants at JEA (Spring) 2019 for their helpful comments. I also thank Aya Tamori (INTAGE Research Inc.), Kimiko Kaneyama, Hidemi Motojima, and Nagisa Yoshioka for their excellent assistance with the surveys and Kazuki Motohashi, Tomohiro Suzuki, Ikuya Yamada, and Kotoko Yanagisawa for their excellent research assistance. This study was conducted as part of the project “Promoting evidence-based policy in Japan,” undertaken at RIETI. This study was also supported by Hitotsubashi University and JSPS KAKENHI, grant number 17K18547.

<sup>†</sup> Graduate School of Economics, Hitotsubashi University. 2-1 Naka, Kunitachi, Tokyo 186-8601, Japan. E-mail: hidefumi.yokoo@r.hit-u.ac.jp

# 1 Introduction

Ethical issues often arise when we run randomized field experiments (Glennester, 2017; Haushofer et al., 2019; Ravallion, 2009). One reason behind these issues is that economists sometimes do not inform the research subjects that they are in an experiment (Levitt and List, 2009). Economists are most likely to acquire informed consent for data collection from the subjects but less likely to explain the experimental design to the subjects (Glennester and Powers, 2016; Teele, 2014). This is quite uncommon or not acceptable for randomized controlled trials (RCTs) in medicine.<sup>1</sup> Possibly due to this practice in economics, implementing partners (e.g., governments and NGOs) raise ethical and reputational concerns about running randomized evaluations and sometimes hesitate to conduct them. Since randomized field experiments can benefit society by their ability to cleanly identify causal impact, excessive concerns would constitute a barrier to making effective policies.

To rigorously evaluate policies while easing the disutility of subjects and the concerns of practitioners, I empirically study the ethical concerns held by potential subjects regarding randomized field experiments. I conduct a series of online surveys on ethical concerns for existing randomized field experiments in the field of economics. As a case study, I select the following six studies: Allcott (2011); Fryer et al. (2015); Hanna et al. (2016); Hosono and Aoyagi (2018); Landry et al. (2006); and Thornton (2008).

In my first survey, comprising approximately 2,000 respondents in Japan, I provided brief explanations on the studies and asked if respondents recognized any ethical issues with them. In my second survey, I focused on two studies among six—the most and least concerning studies from the first survey—and explored a method to ease such ethical concerns by modifying each study. To do so, I applied a randomized online survey experiment (Cruces et al., 2013; Kuziemko et al., 2015) in which approximately 2,000 respondents were randomly

---

<sup>1</sup> Following Favereau (2016), I use the term *randomized field experiment* to refer to an experimental design centered on a random assignment of treatments in the field of economics, while the term *randomized controlled trials (RCTs)* is used to refer to that in medicine throughout this paper.

assigned into three treatment groups and a control group and shown different descriptions. This survey design allows me to estimate the causal impact of changing an attribute of the study—for example, the treatment (from economic incentive to information provision) or research design (from a randomized experiment to a before–after study of an intervention without a control group)—on ethical concerns.

The previous studies on the ethics of randomized field experiments are controversial. For example, both Glennerster and Powers (2016) and Teele (2014) provide a framework for thinking about the ethics of randomized evaluations; however, their arguments are different. Teele (2014) concludes that randomized field experiments differ fundamentally from laboratory experiments or observational studies and require informed consent, the full assessment of the risk of the experiment, and nonexploitative participant selection procedures as minimal steps. Conversely, considering the implementation of programs and other methodologies (e.g., quasi-experimental approaches) as the counterfactual, Glennerster and Powers (2016) conclude that while there are ethical issues specific to randomized evaluations, most of them are not unique to this methodology. Relatedly, List (2008) discusses informed consent associated with natural field experiments and argues that the lack of informed consent seems defensible when the research makes participants better off, benefits society, and confers anonymity and just treatment to all subjects.<sup>2</sup> Note that, these studies discuss the ethics of randomized field experiments from a normative point of view.

While normative analyses on economic methodologies are absolutely important, such debates tend to result in two extreme opinions. Unlike the above studies, which conceptually examine the ethics of the experiments, I conduct an empirical study. Using online surveys, I explore what kind of randomized field experiments are considered by laypeople as involving ethical issues and how researchers can alleviate these concerns by modifying their own research plans. From the results of the positive analyses, I present evidence that contributes to the normative analyses of field experiments. Note that the objective of this paper is to

---

<sup>2</sup> Relatedly, O’Flynn et al. (2016) discuss the definition of ethics in the context of randomized evaluations. Groves Williams (2016) discusses ethics in international development evaluation in general.

improve the methods of experimental studies but not ethically criticize individual papers.

In the first survey, respondents are shown a description of the study and asked the following question: “Do you recognize any ethical issues in this study?” Then, they indicate their concern on a five-point scale. From the results, I find that respondents’ concerns vary among experiments. Relatively few respondents (24%) believe that there is an ethical issue involved in a description that summarizes the work of Fryer et al. (2015), who study the effects of a preschool using the Chicago Heights Early Childhood Center (CHECC) project. In contrast, more than 45% of the respondents recognize that there is an issue in a description that summarizes the work of Landry et al. (2006), who study the impact of a lottery incentive on charitable giving. These results suggest that randomized field experiments are ethically evaluated according to their context, outcomes, treatments, and design.

In the first half of the second survey, which focuses on Fryer et al. (2015), respondents are randomly assigned one out of four slightly different descriptions, which present slightly different experimental designs. Then, they are asked the same question as that in the first survey. I obtain several findings from this survey experiment. First, the response to “There is an ethical issue” significantly increases if parental consent is absent. Second, ethical concern increases if participants are selected at random rather than through self-selection. This result implies that randomization within the self-selected subject pool is more acceptable. These results mean tradeoffs among the Hawthorne effect, specific sample problems (Peters et al., 2016), and ethical issues.

In the second half of the second survey, which focuses on Landry et al. (2006), I find a statistically insignificant impact of changing the research methodology from a randomized field experiment to a before–after study of an intervention without a control group. This result suggests that the internal validity of the analysis can be obtained without increasing ethical issues. Conversely, if the outcome variable of the study is changed from charitable giving to garbage sorting, which both can be considered voluntary public goods provision, then ethical concern significantly decreases. These results imply that as long as the purpose

of the research is not the evaluation of a specific program but rather the testing of a specific theory, then it is possible to alleviate the associated concerns by modifying the research topics.

This paper contributes to three strands of the literature. First, this paper relates to the abovementioned debates on the ethical concerns of randomized field experiments (e.g., Glennerster and Powers, 2016; List, 2008; Ravallion, 2009; Teele, 2014). Unlike these normative analyses, two recent papers have reported the results of surveys on the perceptions of randomized field experiments. Meyer et al. (2019) and Mislavsky et al. (2020) conduct surveys asking about the appropriateness and acceptability, respectively, of the hypothetical scenarios of field experiments. The results of the above two papers are, at first glance, contradictory; Meyer et al. (2019) find that respondents are averse to being involved in experiments, while Mislavsky et al. (2020) find that respondents equally accept an experiment to test a policy and a universal implementation of the same policy. This present paper shares the motivation of the study with the above two papers and presents similar surveys. In addition to the acceptability of experiments, this present paper investigates the causal mechanisms of ethical concerns and methods to alleviate them.

Second, this paper contributes to the nascent literature on the design of experiments that decreases ethical concerns. Duflo et al. (2007) recommend encouragement designs when evaluating programs over which randomizations of the treatment itself are not feasible for ethical reasons. Angrist and Imbens (1991) present an experimental design in which an eligible population is randomly selected, but eligible individuals are allowed to freely choose whether to participate in the program. Narita (2021) develops another experimental design in which subjects with an imaginary budget and personalized clearing price purchase treatment assignment probabilities. While these studies focus on subjects' preferences for *treatments* and propose methods to address ethical issues, this present paper further considers preferences for *studies* including research methodologies, topics, sampling methods, and informed consent. As a result, this paper contributes to the literature by proposing practical

methods to address these concerns.

Third, this paper tangentially relates to a growing set of papers on the nature of preferences for policies. For example, Ambuehl and Ockenfels (2017) study the ethical concerns for increasing incentives to human egg donors, while Elías et al. (2019) study preferences for legalizing payments to kidney donors. Other studies investigate preferences for other policies, such as redistribution (e.g., Cruces et al., 2013; Kuziemko et al., 2015) or nudges (e.g., Hagman et al., 2015; Jung and Mellers, 2016; Sunstein et al., 2018). As in Hagman et al. (2015), the survey in this paper includes a question on social comparison nudges. While the literature on public acceptance of nudges focuses on eliciting preferences for treatments and identifying associated individual characteristics, however, the objective of this paper is different and is as mentioned above. Furthermore, this paper evaluates the impact of changing treatments from economic incentives to informational nudges to identify methods to ease ethical concerns regarding experimental studies.

The paper proceeds as follows. Section 2 presents the motivation behind the surveys conducted in Japan. Section 3 presents the descriptions of the six experiments examined in the first survey and presents the results. Section 4 explains the design of the second survey, describes the data, and presents the main results. Section 5 discusses the implications of the findings, and discusses the limitations of the present study. Section 6 concludes the paper.

## **2 Randomized Field Experiments and Preferences in Japan**

In the last ten years, there has been a rise in the use of randomized field experiments by the Japanese government. First, the Japan International Cooperation Agency (JICA) has started to run randomized evaluations in developing countries. The JICA, mostly in collaboration with economists, has evaluated its programs in various countries, such as Burkina Faso, Cote d’Ivoire, Indonesia, Mongolia, Mozambique, and Niger.<sup>3</sup>

---

<sup>3</sup> One of the early randomized intervention was started in 2010 in Burkina Faso which evaluated the effects of a school-based management (SBM) program (Sawada et al., 2022). Kozuka (2018) also evaluates the SBM program in Niger. Takahashi et al. (2019) evaluate agricultural training in Cote d’Ivoire which is

Second, a series of field experiments were conducted to curb electricity use in Japan. Ito et al. (2018) study one of those experiments that compared the impact of moral suasion and critical peak pricing on electricity demand in Kyoto Province in 2012.<sup>4</sup> The program was designed and jointly implemented by the authors in collaboration with the Ministry of Economy, Trade, and Industry of Japan (METI), a local government, and several private companies. Subsequently, in 2015, the METI evaluated OPOWER’s Home Energy Report (HER) in a northern province in Japan by referring to Allcott (2011).<sup>5</sup>

Third, in 2017, the Japanese government launched a so-called nudge unit, which runs several randomized field experiments (Behavioral Sciences Team, 2019).<sup>6</sup>

All of these movements of evidence-based policy making have brought about an increase in discussions about ethical issues, which arise when we run randomized field experiments among policy makers and researchers (see, for example, Behavioral Sciences Team, 2019). This discussion motivated me to conduct the present study in Japan.

### 3 Survey on Six Experiments

#### 3.1 Survey Data Collection

The first survey was designed by the author and implemented in Japan by the survey company INTAGE Research Inc. in March 2017.<sup>7</sup> In my study, potential respondents in the panel were randomly selected with weights to create a representative Japanese sample in terms of residential area, gender, and age group.

The survey request was sent by email to randomly chosen candidates. I requested the

---

provided in the project by JICA. Tanaka et al. (2019) evaluate the effects of showing leaflets to encourage participation in the public pension system by self-employed workers in Mongolia.

<sup>4</sup> Throughout this paper, I use the term *province* to indicate regions and local governments in Japan although the actual administrative term is *prefecture*. Yamaguchi et al. (2018) argue that *province* is more intuitive for most readers.

<sup>5</sup> Jyukankyo Research Institute Inc. (2016) reports the result of this randomized evaluation.

<sup>6</sup> The Nudge Unit Japan is named the Behavioral Sciences Team (BEST).

<sup>7</sup> INTAGE Holdings which includes INTAGE Research Inc., founded in 1960 in Japan, is ranked 9th in the 2017 American Marketing Association (AMA) Gold Global Top 25 Market Research Firms.

company to implement a sample size of 2,000. In response to this request, the company sent invitation emails to 6,698 candidates. Those who decided to participate were accepted until the number of respondents reached a set number (not known by the author). As a result of this procedure, the sample size for the first survey is 2,107. Prior to the survey, respondents were told that their responses would be used by research institutions, local governments, companies, etc., and they gave their consent.<sup>8</sup>

### 3.2 The Six Randomized Field Experiments Used in the First Survey

[Table 1]

For the first survey, I selected six experimental studies based on the following two criteria: whether the experiments seem to involve ethically sensitive issues and their relevance to current policy discussions in Japan. The selected studies are shown and summarized in Table 1. Half of the selected studies relate to human capital issues: health status, disease testing, or preschools. This reflects that in general, people care more about topics involving human life and death and childhood circumstances. Three experiments are conducted in developed countries, while the other three are conducted in developing countries. This reflects a balance between the increased usage of randomized field experiments in developing countries and the focus of the present study being ethical concerns recognized in a developed country. Finally, a project conducted by the JICA is included as an example for the experiment conducted by Japanese organizations.

In principle, I attempted to summarize the experiments described in original articles as accurately as possible. However, I made several modifications to the original experimental designs, which are mentioned below and in Online Appendix. Most of the modifications were made to simplify the descriptions to make them easy for respondents to understand. In the descriptions, I kept the authors of the six papers anonymous. Throughout the surveys in the present study, I avoided using the word “experiment,” and instead, I used “study” and

---

<sup>8</sup> See Online Appendix A for more details on data collection.



“project.” The Japanese version of the six descriptions was used.

Respondents were shown three randomly assigned descriptions of studies in random order and answered questions for each.<sup>9</sup> For each description, respondents were asked: Do you recognize any ethical issues in this study? Respondents chose one of five options (from “There is a major ethical issue” to “There is no ethical issue at all”). The selected six studies are as follows.

### **Study on a preschool: Fryer et al. (2015)**

The first study I chose is the CHECC project. This project conducts randomized field experiments to evaluate early childhood education interventions. For example, a child and their parents are randomized into one of three groups—preschool treatment, parent academy treatment, and control (Cappelen et al., 2020)<sup>10</sup>—where the preschool used is established for the purpose of this experiment (see Gneezy and List, 2013).<sup>11</sup> Various outcomes are examined, such as cognitive and noncognitive test scores (Fryer et al., 2015), time preferences (Andreoni et al., 2019), risk preferences (Andreoni et al., 2020), and social preferences (Cappelen et al., 2020). I summarized the project and prepared a description of it with several simplifications. The descriptions used in the survey are attached in Online Appendix B. Remarks on the modifications made to the original experiments are shown in Online Appendix C.

For the sampling method of the CHECC project, I described it as “*Parents and children, for a total of 140 families, applied for admission.*” For informed consent, I explicitly mentioned as follows: *Note that the parents of the 140 children who became subjects of the study received an explanation regarding them being the subjects of the study, and they gave their*

---

<sup>9</sup> To keep the time for reading the survey materials and responding to questions short, I provided three randomly assigned descriptions, instead of all the six descriptions, per respondent. Note that three additional descriptions of another study (not shown in this paper) are also shown to each respondent. Furthermore, the respondents were asked two questions for each description. In this paper, survey responses to only one of the two questions is used. In total, respondents were shown six descriptions and answered 12 questions for each. See Section 3.3 for the average duration of the survey.

<sup>10</sup> Fryer et al. (2013) provide an outline of the project, especially in the early stage.

<sup>11</sup> The work of Gneezy and List (2013) was translated into Japanese and published in 2014.

*consent*. For an implementer of the project, I anonymized and framed it as “Professor X.”<sup>12</sup>

Note that, I focused on academic achievement and income, meaning that Fryer et al. (2015) study is the closest among the existing papers produced from the CHECC project. For this reason, to intuitively label the description, I refer to it as “Fryer, Levitt, and List (2015).”<sup>13</sup>

### **Study on HIV testing: Thornton (2008)**

The second study I chose is the study on HIV testing for AIDS prevention in the developing world. Thornton (2008) analyzes the dataset collected in the experiment, which randomly assigned monetary incentives to learn the results of HIV testing. The sample of the study consists of 2,812 individuals in rural Malawi who accepted an HIV test and the followup survey. Thornton (2008) evaluates the impact of incentives on the demand for learning HIV status and subsequent behaviors.<sup>14</sup>

### **Study on charitable giving: Landry et al. (2006)**

The third study I chose is on voluntary contributions to public goods. Landry et al. (2006) conducted a randomized field experiment to study the impact of lotteries on charitable giving. They conducted door-to-door fundraising in North Carolina, where 44 solicitors approached 4,833 households. For households in one among four randomly assigned groups, the single-prize lottery treatment was offered; donors were provided a ticket for a raffle where the winner would receive a USD 1,000 prepaid credit card.<sup>15</sup>

---

<sup>12</sup> For the other five descriptions, however, an implementer of the program is framed differently for randomly selected respondents. More precisely, 75.0% of the descriptions used in the first survey mention the implementer of the program as “Professor X,” but the rest of them purposely mention a different implementer. See Column (9) in Table 1. I intentionally and randomly made this difference to examine another research question. Throughout the paper, I focus on the comparison of six studies and leave the discussion on the impact of the implementer to Online Appendix G. In the regression analysis in this section, I control for this randomness to focus on the comparison of six studies, holding the difference in implementers constant.

<sup>13</sup> Note that, however, Fryer et al. (2015) focus on the parent academy treatment instead of the preschool treatment. As the first description is labeled “Fryer et al. (2015),” five other descriptions are also labeled by the representative papers.

<sup>14</sup> Thornton (2012) and Godlonton and Thornton (2013) use the dataset collected through the same project (the Malawi Diffusion and Ideational Change Project).

<sup>15</sup> Following this study, Landry et al. (2010) conducted another experiment to examine the dynamics of charitable fundraising. Carpenter and Matthews (2017) also study an impact of lotteries on charitable giving. Various other studies conducted randomized field experiment using door-to-door fundraising, for

In the present study, I chose two groups in the original experiment (a voluntary contributions mechanism without seed money and the single-prize lottery) to simplify the description.<sup>16</sup> For the objective of the experiment, I described it as “*to obtain more donations.*” Note that in contrast to the previous two studies, the subjects of Landry et al. (2006) are not informed that such solicitation is part of a research project; thus, it is considered a natural field experiment in the parlance of Harrison and List (2004). I explicitly mentioned this feature in the description as follows: *Note that the 4,800 households that were solicited for donations were not informed of their involvement in the study.*

### **Study on electricity conservation: Allcott (2011)**

The fourth study I chose is on the nudge to encourage electricity conservation. Allcott (2011) evaluates the program that sent Home Energy Report (HER) letters to households. The HER consists of two components: the social comparison module, which compares households’ electricity use to that of their neighbors, and the action steps module, which includes energy conservation tips.<sup>17</sup> I mentioned both of them in the description.

### **Study on household air pollution from cooking: Hanna et al. (2016)**

The fifth study I chose is the evaluation of a program to reduce household air pollution in a developing country. Hanna et al. (2016) evaluate a program implemented in India, where improved cooking stoves are distributed almost for free.<sup>18</sup>

Note that unlike the other five experiments, this program was designed to rollout the treatment, meaning that households in the control group also received stove construction afterward. According to Duflo et al. (2007), such an experimental design, which randomizes the order of phase-in, is considered the fairest way to implement programs.

### **Study on recyclable waste sorting: Hosono and Aoyagi (2018)**

---

example, Soetevent (2011), DellaVigna et al. (2012), and Edwards and List (2014).

<sup>16</sup> The other two treatments are a voluntary contributions mechanism with seed money and the multiple-prize lottery.

<sup>17</sup> Other papers that experimentally evaluate the HER include Ayres et al. (2012), Costa and Kahn (2013), and Allcott and Kessler (2019).

<sup>18</sup> Other papers that study the impact of improved cooking stoves include Mobarak et al. (2012), Bensch and Peters (2015), and Jeuland et al. (2020).

The last study I chose is an experiment implemented by a Japanese organization. Hosono and Aoyagi (2018) analyze a dataset collected in a project conducted by the JICA in Mozambique. In the project, the JICA attempted to encourage household waste-sorting behavior.<sup>19</sup> A total of 1,000 households in a suburb of Maputo are randomly assigned to one of the four groups. Three treatments are evaluated to encourage the sorting of recyclable waste (e.g., plastics and aluminum) from other garbage. In the present study, I focus on in-kind incentive treatment and control groups.

### 3.3 Data and Summary Statistics

[Table 2]

Table 2 shows the characteristics of the sample that completed the first survey. On average, 48% of the respondents are women, 61% are married, 38% live with children, and their average age is 46.7. In addition, I collected information on the time spent on the survey. The median time is 3.4 minutes, and the average is 23 minutes.<sup>20</sup>

### 3.4 Descriptive Results

[Figure 1]

Figure 1 shows the distribution of the responses. Panel A shows the result for Fryer et al. (2015), where approximately 32% of the respondents recognize that there is no ethical issue (Unethical Rating 1 and 2), 44% feel neutral (Unethical Rating 3), and 24% recognize that there is an issue (Unethical Rating 4 and 5). A similar but slightly worse result is obtained for Panel D of Allcott (2011), where approximately 29% recognize that there is no ethical issue, while 27% recognize that there is an issue. For Thornton (2008), the result shows that

---

<sup>19</sup> Chong et al. (2015) evaluate recycling campaigns conducted by an NGO in Peru by using randomized field experiments. Other papers experimentally study interventions to encourage household recycling in developed countries include, for example, Schultz (1999) and Koford et al. (2012).

<sup>20</sup> Figure A2 in the Online Appendix shows a histogram of the time spent on the survey. Table A1 in the Online Appendix reports the result of the regression analysis on the characteristics and time spent on the survey. The time is significantly longer if respondents live with children or if they are part-time employees.

approximately 24% recognize that there is no ethical issue, while 32% recognize that there is an issue, which is quite similar to the results for Hanna et al. (2016) and Hosono and Aoyagi (2018). The study that is recognized as the most unethical is Landry et al. (2006), where approximately 13% of respondents recognize that there is no ethical issue, while more than 45% recognize that there is an issue.

### 3.5 Results from Econometric Analysis

To quantitatively compare the ethical concerns among the six studies, I conduct a regression analysis. In this section, I use a dataset compiled by pooling the responses from the sample of 2,107 respondents. Consider an ordered logit model in the latent variable:

$$y_{ij}^* = \sum_{j=1}^5 \beta_j \cdot EXP_j + x_i' \cdot \gamma + \delta \cdot z_{ij} + \varepsilon_{ij}, \quad (1)$$

where  $y_{ij}^*$  denotes the degree of ethical issues in study  $j$  recognized by respondent  $i$ .  $EXP_j$  is a dummy variable indicating study  $j$ , where  $j = 1, \dots, 5$  represents Fryer et al. (2015) to Hanna et al. (2016), respectively.  $x_i$  represents a vector of characteristics,  $z_{ij}$  represents an order when study  $j$  appears in a survey of respondent  $i$ , and  $\varepsilon_{ij}$  is the error term, which is assumed to follow a standard logistic distribution.<sup>21</sup> The five studies are compared to Hosono and Aoyagi (2018) by estimating  $\beta_j$ .

The observed, ordered dependent variable is linked to the latent variable  $y_{ij}^*$  through cut points  $\mu_k$  in the following way:

$$y_{ij} = \begin{cases} 1 & \text{if } y_{ij}^* < \mu_1, \\ k & \text{if } \mu_{k-1} \leq y_{ij}^* < \mu_k \text{ where } k = \{2, 3, 4\}, \\ 5 & \text{if } \mu_4 \leq y_{ij}^*. \end{cases}$$

Columns (1) and (2) in Table 3 report the estimation results. The estimated coefficients

---

<sup>21</sup> In the first survey, the variable *Order* is an integer ranging from one to six. See Footnote 9 for more information.

for Fryer et al. (2015) and Allcott (2011) are negative and significant, meaning that respondents on average recognize less ethical issues in them compared to Hosono and Aoyagi (2018). Landry et al. (2006) is significantly positive, meaning that significantly large ethical issues are recognized. The coefficients for Thornton (2008) and Hanna et al. (2016) are close to zero and not significant at the 10% level, meaning that ethical issues are almost similar to those of Hosono and Aoyagi (2018). *Order* is statistically significantly negative, meaning that the recognition of ethical issues is small for the same experiment if it is shown later in the survey.

[Table 3]

Columns (3) and (4) in Table 3 report the results from linear regressions of Equation (1). The signs and statistical significance of the coefficients are similar to the ordered logit results. The constant term in Column (3) is 3.2, meaning that Hosono and Aoyagi (2018) is, on average, recognized as “3: Neutral” or slightly worse. Since the coefficient for Fryer et al. (2015) is  $-0.22$ , its average ethical issue is 2.8 on a five-point scale. The coefficient for Landry et al. (2006) is 0.37 and significant. Columns (2) and (4) consistently show that women are more likely to recognize ethical issues than are men. Age is positively associated with ethical concerns. The coefficient for *Order* is  $-0.03$ .

## 4 Randomized Survey Experiments on Two Experiments

### 4.1 Overview of the Second Survey

The results of the previous section show that Fryer et al. (2015) is recognized as having the least ethical issues, while Landry et al. (2006) is recognized as having the most ethical issues among the six studies. Why do ethical concerns vary among the experiments? Can we alleviate these concerns by modifying the original studies?

To investigate the above questions, the second survey was designed. It was implemented in March 2018 by INTAGE Research Inc. In the second survey, I focus on the two studies

as a contrasting example and use the design of a randomized online survey experiment. I develop three hypotheses, as described below, for each study. Using the same procedure as that of the first survey, 2,146 respondents are invited to take the second survey and are randomly shown one of four descriptions in each study.<sup>22</sup>

## 4.2 Hypotheses and Treatments

### 4.2.1. Hypotheses About Small Ethical Concerns in Fryer et al. (2005)

Examining the description used in the first survey leads us to several hypotheses regarding why the work of Fryer et al. (2015) is ethically more acceptable than are other studies. First, informed consent may matter. As mentioned in Section 3.2, the last sentence in the description explicitly mentions informed consent in this experiment. The presence of consent from subjects, which is missing in Landry et al. (2006), could have alleviated the recognition of an ethical issue. The first treatment tests this hypothesis by deleting this last sentence from the description.

Second, the sampling strategy may matter. In the CHECC project, households were recruited to the project, applied according to their decision, and were randomly assigned to control and treatment groups (Fryer et al., 2015; Cappelen et al., 2020). In contrast, the samples in Landry et al. (2006) did not request to be solicited but became targets of door-to-door fundraising. Thus, in the second treatment, I modified the sentence to mention that subjects were defined by a researcher rather than by applicants as subjects: *parents and their children from 140 families living in the area are defined as the research subjects.*

Third, the existence of a followup for the control group may matter. In the CHECC project, parents and their children in the control group are also invited to holiday parties (see, Gneezy and List, 2013). This may be recognized as compensation to the control group and alleviate the issues. In the third treatment, I deleted the sentences on invitations to holiday parties.

---

<sup>22</sup> Respondents for the first survey are intentionally excluded from the second survey. The second online survey was conducted from March 2 to March 5, 2018.

#### 4.2.2. Hypotheses About Large Ethical Concerns in Landry et al. (2006)

To investigate methods to ease ethical concerns by modifying Landry et al. (2006), I developed three hypotheses. First, the research design of a randomized field experiment may increase the recognition of ethical issues. To test whether it is worse than other research designs in terms of ethical concerns, I changed the program evaluation methodology to a before–after study of an intervention without a control group.

Second, the treatment may matter. Landry et al. (2006) use a raffle to encourage donations. As previous studies discuss a crowding out of intrinsic altruism by extrinsic incentives (e.g., Bénabou and Tirole, 2006), people may not like this treatment as a means of fostering prosocial behavior. Alternatively, in the second treatment, I change the treatment to social comparison information that is used in Allcott (2011) and others for energy conservation and Frey and Meier (2004) and Shang and Croson (2009) for charitable giving. Specifically, I mention that donations are collected with flyers where a message of “*In the neighboring town, 80% of the households donated*” is printed.

Third, the topic of the study may matter. Studies to encourage charitable giving may be recognized as unethical, regardless of how we encourage or evaluate such programs. Theoretically, charitable giving is modeled as the private provision of public goods (e.g., Bergstrom et al., 1986). Similarly, household waste sorting to decrease social cost to the environment is also modeled as the private provision of public goods (e.g., Brekke et al., 2003). Therefore, in the third treatment, I change the topic of study from charitable giving to waste sorting, which is similar to Hosono and Aoyagi (2018). Note that I keep the treatment and methodology of program evaluation unchanged. More specifically, I mention that Professor X collaborates with a city government and calls for sorting food waste from other garbage. In the campaign, households are “*asked to sort with a raffle in which one among all recyclers could win JPY 100,000.*”

Based on the above hypotheses, I modified the description used in the first survey. As a



result, I prepared four descriptions for each study.<sup>23</sup> Respondents in the control group are shown the same descriptions with the first survey. Respondents are randomly assigned to one group among four groups for each study. This results in 2,146 respondents being randomly assigned to 16 groups. As in the first survey, the orders in which Fryer et al. (2015) and Landry et al. (2006) have been shown are randomly determined.<sup>24</sup>

### 4.3 Verifying Randomizations

Tables A2 and A3 in Online Appendix present summary statistics for respondents of randomized survey experiments on Fryer et al. (2015) and Landry et al. (2006), respectively. The results show that only three and four differ at  $p < 0.10$  out of 39 differences each in Fryer et al. (2015) and Landry et al. (2006), respectively. From these, I conclude that the four groups in each survey experiment are very similar.

### 4.4 Main Results

To evaluate the causal impacts of the treatments, I estimate the models of ordered logit and OLS separately for the samples in each study:

$$y_{ij}^* = \beta_1 \cdot T_1 + \beta_2 \cdot T_2 + \beta_3 \cdot T_3 + \delta \cdot z_{ij} + \varepsilon_i, \quad \text{for } j = 1 \text{ or } 3, \quad (2)$$

where  $T_n$  is a dummy variable indicating treatment  $n$ .<sup>25</sup>

[Table 4]

Table 4 reports the results of the survey on Fryer et al. (2015). I compute  $p$ -values based on the randomization inference procedure of Young (2019) for individual treatment

---

<sup>23</sup> These descriptions are attached in Online Appendix D.

<sup>24</sup> In this second survey, I use “Professor X” as the implementer of the program for all the descriptions.

<sup>25</sup> Figures A3 and A4 in Online Appendix show the distribution of the responses to the second survey on Fryer et al. (2015) and Landry et al. (2006), respectively. Interestingly, the distribution of the responses to the waste-sorting version of Landry et al. (2006)(Figure A4 Panel D) is closer to that of Hosono and Aoyagi (2018)(Figure 1 Panel F) rather than that of Landry et al. (2006) in the first survey (Figure 1 Panel C).

effects. I also report the results adjusting for multiple hypothesis testing using the procedure of Westfall and Young (1993) under the null hypothesis that all treatment effects in the equation are zero.

The results show that deleting the sentence on informed consent by parents increases ethical concerns (the randomization- $t$   $p$ -value of 0.007, column 1). The coefficient for this treatment is 0.17 for the OLS (column 3). This magnitude of the effect is similar to the difference between Fryer et al. (2015) and Thornton (2008) in the first survey (Table 3, column 3). If the sample is selected by a researcher, irrelevant to one’s willingness to participate, then ethical concerns increase ( $p$ -value 0.016), while the magnitude of the effect is slightly smaller than that of treatment 1. Deleting the sentence on holiday parties in which control groups are also invited does not increase these concerns ( $p$ -value 0.963). From the results, adjusting for multiple hypothesis testing, I can reject the null hypothesis that all treatment effects are zero ( $p$ -value 0.018). Finally, *Order* negatively affects the recognition of ethical issues, which is consistent with the first survey.

[Table 5]

Table 5 reports the results of the survey on Landry et al. (2006). Changing the methodology of program evaluation from a randomized field experiment to a before–after study does not decrease ethical concerns (the randomization- $t$   $p$ -value of 0.478, column 1). Changing treatments from a raffle to social comparison message slightly and weakly decreases ethical concerns ( $p$ -value 0.097). Finally, changing a topic of the study from encouraging charitable giving to waste sorting decreases ethical concerns ( $p$ -value 0.000 for the individual coefficient and 0.001 for the result, adjusting for the Westfall-Young multiple testing). The magnitude of this effect is large. The coefficient for this treatment is  $-0.21$  for the OLS (column 3), which accounts for more than half of the difference between Landry et al. (2006) and Hosono and Aoyagi (2018) in the first survey (the coefficient of 0.37, Table 3, column 3).

## 4.5 Subgroup Analysis

[Table 6]

Table 6 reports the results of subgroup analyses to examine whether there is heterogeneity in impacts by gender. The top panel reports that women are significantly affected by the treatments in Fryer et al. (2015).<sup>26</sup> The result for men shows no significant impacts for any of the three treatments. Moreover, men are not affected by the order of the survey, while women are affected significantly.

The bottom panel reports that, unlike in Table 5, changing a raffle to a message in Landry et al. (2006) significantly decreases the ethical concerns of women. Furthermore, the magnitude of the effect is not small, as the coefficient for the treatment is  $-0.18$  for the OLS (column 3). For the treatment that changes the topic from charitable giving to waste sorting, the result for women shows significant negative impacts, while men show weakly significant and nonnegligible negative impacts. Overall, there is heterogeneity in the impacts—women are more sensitive than men to the modifications of the studies in terms of ethical concerns.

## 5 Discussion

### 5.1 Robustness Checks Related to the Time Spent on the Surveys

Some respondents may not carefully read the descriptions. I conduct a comparison of the six studies the same way as I did in Section 3 but drop respondents whose time spent on the survey is in the bottom 10% (see Online Appendix Table A4). The result is consistent with Table 3. Moreover, the absolute values of the estimated coefficients are larger than those in Table 3, indicating that differences in ethical concerns among studies become larger if we focus on respondents who take a long time to complete the survey.

---

<sup>26</sup> I also report  $p$ -values adjusted for multiple-hypothesis testing using the procedure of Westfall and Young (1993) and Young (2019) within two regressions of a same model for women and men (e.g., columns 1 and 4 in the top panel). I can reject the null hypothesis that all treatment effects are zero for both men and women ( $p$ -value of 0.001 for columns 1 and 4). Similarly, I can reject the null hypothesis for the bottom panel as well ( $p$ -value 0.003).

Similarly, I analyze the two randomized survey experiments considering the time spent on the survey. For the dataset used in Section 4, I create a dummy variable that takes a value of one if time spent on the survey is longer than the median and zero otherwise (*Long time*). Table A5 in the Online Appendix shows the results of the analyses incorporating the interaction terms of treatments and *Long time*. For Fryer et al. (2015), the interaction terms are consistent with Table 4, while treatments without interaction with *Long time* are not significant. This result can be interpreted as those who read the description carefully being more sensitive to the lack of informed consent or self-selection into the experiment.<sup>27</sup>

The result is slightly different for Landry et al. (2006). Table A6 shows the result, which is consistent with Table 5 for treatment 3 (waste sorting rather than donations) *without* the interaction. This suggests that those who read the description quickly find fewer ethical issues when glancing a waste sorting study; however, changing the topic is not enough to alleviate the concerns of those who read the description carefully. Finally, changing the design of the study to a before–after study does not affect the concerns within each of the two groups (*Long time* = 0 or 1), suggesting no heterogeneous effects and an average effect.

## 5.2 Interpretations and Implications of the Results

Several implications are obtained from a series of surveys. From the first survey, I find that the distribution of the recognition of ethical issues widely varies among the six studies. Implementing partners frequently raise ethical concerns about randomized evaluations in general; however, whether subjects identify ethical issues depends on the experiments. Not all field experiments but some specific topics, treatments and designs involve ethical issues. In a specific worst case, researchers are required to modify research plans to improve social welfare through research activities.

At first glance, experiments that may affect lifetime success, such as early childhood

---

<sup>27</sup> Interestingly, the coefficient for *Long time* is negative and significant. This correlation can be interpreted in two ways. First, those who recognize relatively large ethical issues are more likely to quickly read through the description. Second, those who spend a longer time reading the description recognize less ethical issues as a result of reading it carefully.

interventions, seem ethically contentious. However, the results reveal that the number of respondents who recognize ethical issues is the lowest for the CHECC project. Respondents may balance the risks and benefits of the experiment considering whether the findings from the experiment are beneficial and relevant to their lives. Another explanation is that a situation where only half of the applicants are admitted to attend a preschool is common and unsurprising for the respondents since the demand for subsidized childcare often exceeds supply in Japan (for more details, see Yamaguchi et al., 2018). People may be more likely to accept an experiment if the partial and random assignment of the treatment is a familiar situation for the context and culture of their lives.

The second survey partly identifies the reasons for low ethical concerns in the CHECC project. First, women recognize more ethical issues if there is no sentence on informed consent. Among the six studies, the CHECC is the only experiment that informed the subjects of the objective of the study and acquired consent. Note that for the descriptions of the other five experiments, I explicitly mentioned that the subjects were not fully informed of the objectives and designs of the studies (see Online Appendix B). This result empirically supports the normative discussions in the literature on the importance of informed consent (e.g., Glennerster and Powers, 2016; Teele, 2014). Second, respondents (especially women) recognize fewer ethical issues if subjects voluntarily participate in an experiment based on their decisions compared to researchers randomly selected from the population. Taken together, the random assignment of treatments over subjects who agreed to be in the experiment is recognized as being better from an ethical perspective.

These findings pose tradeoffs to randomized field experiments. Informing subjects that they are taking part in an experiment may change their behavior (Duflo et al., 2007; Harrison and List, 2004). So-called Hawthorne and John Henry effects can occur when we acquire informed consent and may limit the external validity of experiments. Similarly, self-selection into an experiment often makes the sample different from the policy population, which results in biases in the estimate and limits external validity (Deaton, 2010; Peters et al.,

2016). Apparently, researchers and implementers face a difficult problem of balancing the external validity of the result and ethical concerns of subjects when running randomized evaluations.

Among the six examined studies, Landry et al. (2006) is recognized as the least acceptable from an ethical perspective. The result of the second survey suggests that respondents do not recognize concerns because the researcher randomizes the treatment. Possibly, however, respondents are concerned with the research question itself; that is, “Can we encourage charitable giving by a raffle?” One interpretation of the result is that respondents believe that it is unethical to incentivize charitable giving. My result shows that it is less problematic if subjects are solicited using a message with a nudge. Previous studies examine the crowding-out of intrinsic motivations to donate by monetary incentives (e.g., Mellström and Johannesson, 2008). People may dislike being incentivized to make donations.

This result implies that the preferences for experiments are associated with the preferences for treatments. The result also suggests that preferences are associated with the type of outcome variables. Holding the treatment constant and changing the outcome variable from charitable giving to waste sorting cease such concerns. This suggests that if the motivation to use the experiment is not an evaluation of a program (e.g., a raffle to encourage charitable giving) but rather a test of a theory (e.g., a model of voluntary provision of public goods), then we can alleviate ethical concerns by changing the topic and context of the study.

### 5.3 Study Limitations

Some limitations in the present paper are worth noting. First, I compare only six randomized field experiments among thousands implemented or analyzed in the field of economics.<sup>28</sup> Second, while my randomized survey experiments partly unmask the reasons for relatively low or high ethical concerns for specific studies, the findings in the present study cannot fully

---

<sup>28</sup> There are 7,030 studies registered in the AEA RCT Registry as of April 12, 2023. Peters et al. (2016) review 92 papers that used a randomized field experiment and were published between 2009 and 2014. Lewis and Rao (2015) review 25 randomized field experiments that measure returns to digital advertising.

explain the large difference between Fryer et al. (2015) and Landry et al. (2006). Relatedly, I show that the examined strategies can alleviate concerns; however, the magnitude of such alleviation is not large. Third, there may exist disutility other than that represented by ethical concerns. For example, subjects may feel anxiety due to being treated by untested treatments. Furthermore, subjects may find disutility from the inequality of treatment status as a result of a random assignment. These types of possible disutility are not examined in the present paper.

## 6 Conclusions

Randomized field experiments can improve social welfare by rigorously evaluating policies or testing economic theories. However, there is a concern that experiments may generate utility loss for subjects and implementing partners. In this study, I conduct an online survey to compare potential subjects' ethical concerns with six previous experiments in the field of economics. I find that the degree of ethical concerns varies among respondents and experiments. A certain proportion of respondents are very concerned, while others are not. Both researchers and practitioners need to take into account this heterogeneity in preferences for economic studies.

From two randomized survey experiments, I find that it is possible to alleviate concerns by modifying research projects. However, the strategies to alleviate the concerns bring about tradeoffs. Easing ethical concerns results in decreasing the external validity of the randomized evaluation design. Note that, the method of this paper can be used to understand the utility or disutility of field experiments for citizens not only *ex post* but also *ex ante*. Future tasks include conducting similar surveys in other countries and examining other experiments both before and after the interventions.

As emphasized by Glennerster and Powers (2016), balancing the risks and benefits of research is required for economists to improve social welfare through their experiments. This paper reveals that randomized experiments are useful for examining a wide range of

issues, including the ethical issues involved in this method. Thus, we need to improve this beneficial method and reduce the risks involved in it to further utilize it.

## References

- Allcott, Hunt (2011). “Social norms and energy conservation.” *Journal of Public Economics*, 95(9-10), 1082–1095.
- Allcott, Hunt and Judd B. Kessler (2019). “The Welfare Effects of Nudges: A Case Study of Energy Use Social Comparisons.” *American Economic Journal: Applied Economics*, 11(1), 236–76.
- Ambuehl, Sandro and Axel Ockenfels (2017). “The Ethics of Incentivizing the Uninformed: A Vignette Study.” *American Economic Review: Papers & Proceedings*, 107(5), 91–95.
- Andreoni, James, Amalia Di Girolamo, John A. List, Claire Mackevicius, and Anya Samek (2020). “Risk preferences of children and adolescents in relation to gender, cognitive skills, soft skills, and executive functions.” *Journal of Economic Behavior & Organization*, 179, 729–742.
- Andreoni, James, Michael A. Kuhn, John A. List, Anya Samek, Kevin Sokal, and Charles Sprenger (2019). “Toward an understanding of the development of time preferences: Evidence from field experiments.” *Journal of Public Economics*, 177, 104039.
- Angrist, Joshua D and Guido W Imbens (1991). “Sources of identifying information in evaluation models.” Technical Working Paper 117, National Bureau of Economic Research. Available at <https://www.nber.org/papers/t0117>.
- Ayres, Ian, Sophie Raseman, and Alice Shih (2012). “Evidence from Two Large Field Experiments that Peer Comparison Feedback Can Reduce Residential Energy Usage.” *Journal of Law, Economics, and Organization*, 29(5), 992–1022.
- Behavioral Sciences Team (2019). “Annual Report (FY2017 and FY2018).” Report, Behavioral Sciences Team (BEST), Government of Japan. Available at [http://www.env.go.jp/earth/ondanka/nudge/report1\\_Eng.pdf](http://www.env.go.jp/earth/ondanka/nudge/report1_Eng.pdf).
- Bénabou, Roland and Jean Tirole (2006). “Incentives and Prosocial Behavior.” *American Economic Review*, 96(5), 1652–1678.
- Bensch, Gunther and Jörg Peters (2015). “The intensive margin of technology adoption—Experimental evidence on improved cooking stoves in rural Senegal.” *Journal of Health Economics*, 42, 44–63.
- Bergstrom, Theodore, Lawrence Blume, and Hal Varian (1986). “On the private provision of public goods.” *Journal of Public Economics*, 29(1), 25–49.



- Brekke, Kjell Arne, Snorre Kverndokk, and Karine Nyborg (2003). “An economic model of moral motivation.” *Journal of Public Economics*, 87(9), 1967–1983.
- Cappelen, Alexander, John List, Anya Samek, and Bertil Tungodden (2020). “The Effect of Early-Childhood Education on Social Preferences.” *Journal of Political Economy*, 128(7), 2739–2758.
- Carpenter, Jeffrey and Peter Hans Matthews (2017). “Using raffles to fund public goods: Lessons from a field experiment.” *Journal of Public Economics*, 150, 30–38.
- Chong, Alberto, Dean Karlan, Jeremy Shapiro, and Jonathan Zinman (2015). “(Ineffective) Messages to Encourage Recycling: Evidence from a Randomized Evaluation in Peru.” *World Bank Economic Review*, 29(1), 180–206.
- Costa, Dora L and Matthew E Kahn (2013). “Energy conservation “nudges” and environmentalist ideology: Evidence from a randomized residential electricity field experiment.” *Journal of the European Economic Association*, 11(3), 680–702.
- Cruces, Guillermo, Ricardo Perez-Truglia, and Martin Tetaz (2013). “Biased perceptions of income distribution and preferences for redistribution: Evidence from a survey experiment.” *Journal of Public Economics*, 98, 100–112.
- Deaton, Angus (2010). “Instruments, Randomization, and Learning about Development.” *Journal of Economic Literature*, 48(2), 424–455.
- DellaVigna, Stefano, John A. List, and Ulrike Malmendier (2012). “Testing for Altruism and Social Pressure in Charitable Giving.” *Quarterly Journal of Economics*, 127(1), 1–56.
- Duflo, Esther, Rachel Glennerster, and Michael Kremer (2007). “Using Randomization in Development Economics Research: A Toolkit.” In *Handbook of Development Economics*, vol. 4, edited by T. Paul Schultz and John A. Strauss, pp. 3895–3962. Elsevier.
- Edwards, James T. and John A. List (2014). “Toward an understanding of why suggestions work in charitable fundraising: Theory and evidence from a natural field experiment.” *Journal of Public Economics*, 114, 1–13.
- Elías, Julio J., Nicola Lacetera, and Mario Macis (2019). “Paying for Kidneys? A Randomized Survey and Choice Experiment.” *American Economic Review*, 109(8), 2855–88.
- Favereau, Judith (2016). “On the analogy between field experiments in economics and clinical trials in medicine.” *Journal of Economic Methodology*, 23(2), 203–222.
- Frey, Bruno S and Stephan Meier (2004). “Social comparisons and pro-social behavior: Testing “conditional cooperation” in a field experiment.” *American Economic Review*, 94(5), 1717–1722.
- Fryer, Roland, Steven Levitt, John List, and Anya Samak (2013). “Chicago Heights Early Childhood Center: Early Results from a Field Experiment on the Temporal Allocation of Schooling.” Available at <https://cpb-us-w2.wpmucdn.com/voices.uchicago.edu/>

dist/f/1276/files/2018/10/CHECC-Presentation-13pmhkk.pdf.

- Fryer, Roland G, Steven D Levitt, and John A List (2015). “Parental incentives and early childhood achievement: A field experiment in Chicago heights.” NBER Working Paper Series 21477, National Bureau of Economic Research. Available at <http://www.nber.org/papers/w21477>.
- Glennerster, Rachel (2017). “The practicalities of running randomized evaluations: Partnerships, measurement, ethics, and transparency.” In *Handbook of Economic Field Experiments*, vol. 1, pp. 175–243. Elsevier.
- Glennerster, Rachel and Shawn Powers (2016). “Balancing risk and benefit: Ethical tradeoffs in running randomized evaluations.” In *Oxford Handbook of Professional Economic Ethics*, edited by George DeMartino and Deirdre McCloskey, pp. 367–401. Oxford University Press, Oxford.
- Gneezy, Uri and John A. List (2013). *The Why Axis: Hidden Motives and the Undiscovered Economics of Everyday Life*. PublicAffairs.
- Godlonton, Susan and Rebecca L. Thornton (2013). “Learning from Others’ HIV Testing: Updating Beliefs and Responding to Risk.” *American Economic Review: Papers & Proceedings*, 103(3), 439–44.
- Groves Williams, Leslie (2016). “Ethics in international development evaluation and research: What is the problem, why does it matter and what can we do about it?” *Journal of Development Effectiveness*, 8(4), 535–552.
- Hagman, William, David Andersson, Daniel Västfjäll, and Gustav Tinghög (2015). “Public Views on Policies Involving Nudges.” *Review of Philosophy and Psychology*, 6(3), 439–453.
- Hanna, Rema, Esther Duflo, and Michael Greenstone (2016). “Up in smoke: The influence of household behavior on the long-run impact of improved cooking stoves.” *American Economic Journal: Economic Policy*, 8(1), 80–114.
- Harrison, Glenn W. and John A. List (2004). “Field Experiments.” *Journal of Economic Literature*, 42(4), 1009–1055.
- Haushofer, Johannes, Michala Iben Riis-Vestergaard, and Jeremy Shapiro (2019). “Is there a social cost of randomization?” *Social Choice and Welfare*, 52(4), 709–739.
- Hosono, Tomoyuki and Keitaro Aoyagi (2018). “Effectiveness of interventions to induce waste segregation by households: Evidence from a randomized controlled trial in Mozambique.” *Journal of Material Cycles and Waste Management*, 20(2), 1143–1153.
- Ito, Koichiro, Takanori Ida, and Makoto Tanaka (2018). “Moral suasion and economic incentives: Field experimental evidence from energy demand.” *American Economic Journal: Economic Policy*, 10(1), 240–67.
- Jeuland, Marc, Subhrendu K. Pattanayak, Jie-Sheng Tan Soo, and Faraz Usmani (2020).

- “Preferences and the Effectiveness of Behavior-Change Interventions: Evidence from Adoption of Improved Cookstoves in India.” *Journal of the Association of Environmental and Resource Economists*, 7(2), 305–343.
- Jung, Janice Y and Barbara A Mellers (2016). “American attitudes toward nudges.” *Judgment & Decision Making*, 11(1), 62–74.
- Jyukankyo Research Institute Inc. (2016). “Improvement of infrastructure to promote rationalization of energy use in fiscal year 2015 (Survey on the effect of providing information on energy use on the promotion of changes in household energy conservation behavior).” Tech. rep., Minister of Economy, Trade and Industry (METI), Japan.
- Koford, Brandon C, Glenn C Blomquist, David M Hardesty, and Kenneth R Troske (2012). “Estimating consumer willingness to supply and willingness to pay for curbside recycling.” *Land Economics*, 88(4), 745–763.
- Kozuka, Eiji (2018). “Enlightening Communities and Parents for Improving Student Learning Evidence from Randomized Experiment in Niger.” JICA-RI Working Paper No. 166, JICA Research Institute.
- Kuziemko, Ilyana, Michael I Norton, Emmanuel Saez, and Stefanie Stantcheva (2015). “How elastic are preferences for redistribution? Evidence from randomized survey experiments.” *American Economic Review*, 105(4), 1478–1508.
- Landry, Craig E, Andreas Lange, John A List, Michael K Price, and Nicholas G Rupp (2006). “Toward an understanding of the economics of charity: Evidence from a field experiment.” *Quarterly Journal of Economics*, 121(2), 747–782.
- Landry, Craig E., Andreas Lange, John A. List, Michael K. Price, and Nicholas G. Rupp (2010). “Is a Donor in Hand Better Than Two in the Bush? Evidence from a Natural Field Experiment.” *American Economic Review*, 100(3), 958–83.
- Levitt, Steven D and John A List (2009). “Field experiments in economics: The past, the present, and the future.” *European Economic Review*, 53(1), 1–18.
- Lewis, Randall A. and Justin M. Rao (2015). “The Unfavorable Economics of Measuring the Returns to Advertising.” *Quarterly Journal of Economics*, 130(4), 1941–1973.
- List, John A (2008). “Informed consent in social science.” *Science*, 322(5902), 672.
- Mellström, Carl and Magnus Johannesson (2008). “Crowding Out in Blood Donation: Was Titmuss Right?” *Journal of the European Economic Association*, 6(4), 845–863.
- Meyer, Michelle N, Patrick R Heck, Geoffrey S Holtzman, Stephen M Anderson, William Cai, Duncan J Watts, and Christopher F Chabris (2019). “Objecting to experiments that compare two unobjectionable policies or treatments.” *Proceedings of the National Academy of Sciences*, 116(22), 10723–10728.
- Mislavsky, Robert, Berkeley Dietvorst, and Uri Simonsohn (2020). “Critical Condition: Peo-

- ple Don't Dislike a Corporate Experiment More Than They Dislike Its Worst Condition." *Marketing Science*, 39(6), 1092–1104.
- Mobarak, Ahmed Mushfiq, Puneet Dwivedi, Robert Bailis, Lynn Hildemann, and Grant Miller (2012). "Low demand for nontraditional cookstove technologies." *Proceedings of the National Academy of Sciences*, 109(27), 10815–10820.
- Narita, Yusuke (2021). "Incorporating ethics and welfare into randomized experiments." *Proceedings of the National Academy of Sciences*, 118(1), e2008740118.
- O'Flynn, Peter, Chris Barnett, and Laura Camfield (2016). "Assessing contrasting strategies for ensuring ethical practice within evaluation: Institutional review boards and professionalisation." *Journal of Development Effectiveness*, 8(4), 561–568.
- Peters, Jörg, Jörg Langbein, and Gareth Roberts (2016). "Policy evaluation, randomized controlled trials, and external validity—A systematic review." *Economics Letters*, 147, 51–54.
- Ravallion, Martin (2009). "Evaluation in the Practice of Development." *World Bank Research Observer*, 24(1), 29–53.
- Sawada, Yasuyuki, Takeshi Aida, Andrew S. Griffen, Eiji Kozuka, Haruko Noguchi, and Yasuyuki Todo (2022). "Democratic institutions and social capital: Experimental evidence on school-based management from a developing country." *Journal of Economic Behavior and Organization*, 198, 267–279.
- Schultz, P. Wesley (1999). "Changing Behavior With Normative Feedback Interventions: A Field Experiment on Curbside Recycling." *Basic and Applied Social Psychology*, 21(1), 25–36.
- Shang, Jen and Rachel Croson (2009). "A Field Experiment in Charitable Contribution: The Impact of Social Information on the Voluntary Provision of Public Goods." *Economic Journal*, 119(540), 1422–1439.
- Soetevent, Adriaan R. (2011). "Payment Choice, Image Motivation and Contributions to Charity: Evidence from a Field Experiment." *American Economic Journal: Economic Policy*, 3(1), 180–205.
- Sunstein, Cass R, Lucia A Reisch, and Julius Rauber (2018). "A worldwide consensus on nudging? Not quite, but almost." *Regulation & Governance*, 12(1), 3–22.
- Takahashi, Kazushi, Yukichi Mano, and Keijiro Otsuka (2019). "Learning from experts and peer farmers about rice production: Experimental evidence from Cote d'Ivoire." *World Development*, 122, 157–169.
- Tanaka, Tomoaki, Junichi Yamasaki, Yasuyuki Sawada, and Khaliun Dovchinsuren (2019). "Barriers to Public Pension Program Participation in a Developing Country." JICA-RI Working Paper No. 199, JICA Research Institute.

- Teele, Dawn Langan (2014). “Reflections on the ethics of field experiments.” In *Field experiments and their critics: Essays on the uses and abuses of experimentation in the social sciences*, edited by Dawn Langan Teele, pp. 115–140. Yale University Press, New Haven & London.
- Thornton, Rebecca L (2008). “The demand for, and impact of, learning HIV status.” *American Economic Review*, 98(5), 1829–63.
- Thornton, Rebecca L. (2012). “HIV testing, subjective beliefs and economic behavior.” *Journal of Development Economics*, 99(2), 300–313.
- Westfall, Peter H and S Stanley Young (1993). *Resampling-based multiple testing: Examples and methods for p-value adjustment*. John Wiley & Sons, New York.
- Yamaguchi, Shintaro, Yukiko Asai, and Ryo Kambayashi (2018). “How does early childcare enrollment affect children, parents, and their interactions?” *Labour Economics*, 55, 56–71.
- Young, Alwyn (2019). “Channeling Fisher: Randomization Tests and the Statistical Insignificance of Seemingly Significant Experimental Results.” *Quarterly Journal of Economics*, 134(2), 557–598.

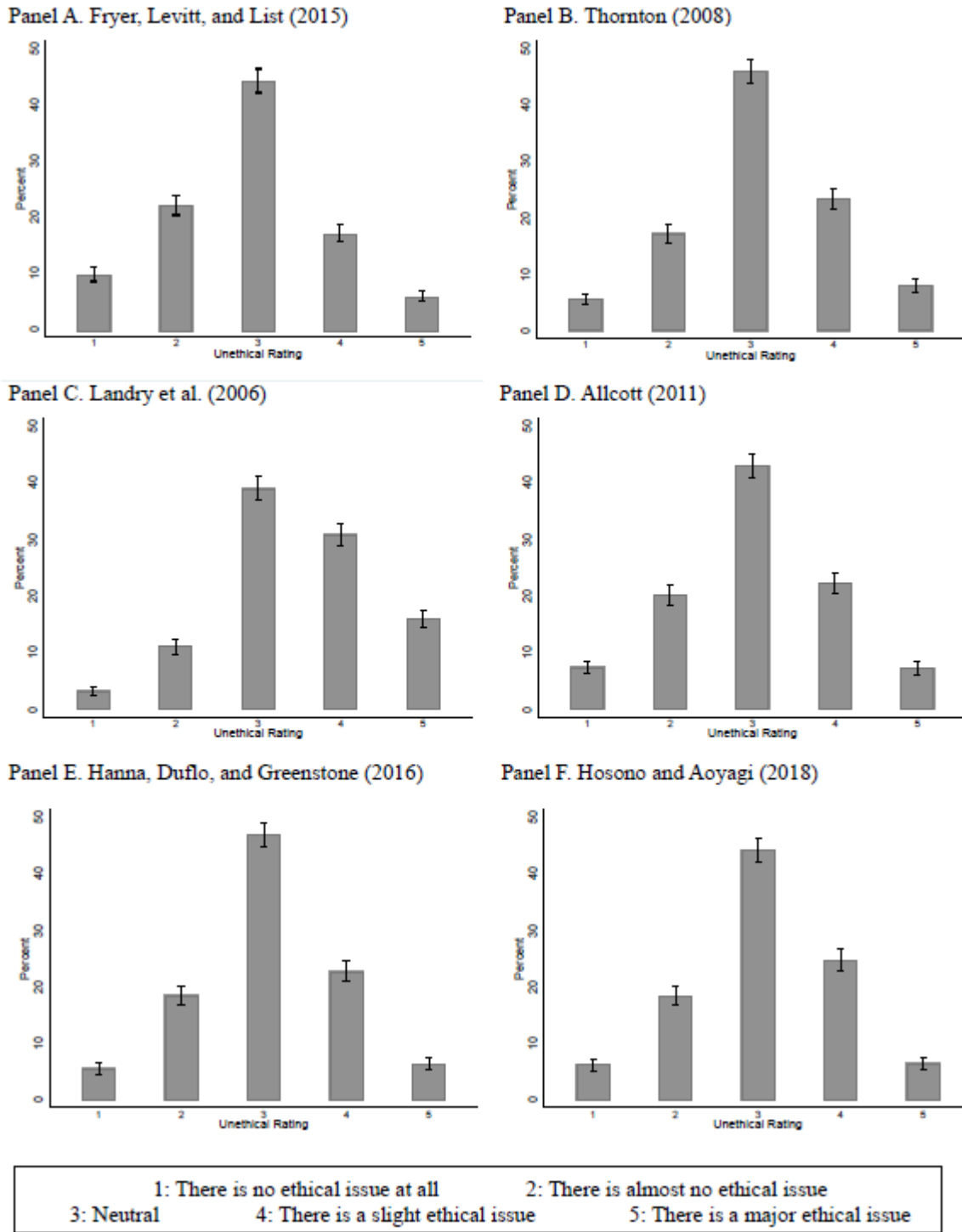


Figure 1: Response to the first survey

*Notes:* This figure shows the distribution (percentages) of the survey response to the question “Do you recognize any ethical issues in this study?” The vertical bars and caps are 95 % confidence intervals. The average number of observations for the six questions is 1,053.5.

Table 1: Summary of the six experiments examined in the first survey

	(1) Label of the descriptions	(2) Outcome variables	(3) Treatments	(4) Sample size	(5) Informed	(6) Monetary incentive	(7) Human capital	(8) Developing countries	(9) Implementer
1	Fryer, Levitt, and List (2015)	Academic achievement and lifetime earnings	Preschool	140	No	No	No	No	Professor X
2	Thornton (2008)	Going to HIV testing centers to be informed	Reward	3,000	No	Yes	Yes	Yes	Professor X
3	Landry et al. (2006)	Donation	Raffle	4,800	Yes	Yes	No	No	Professor X
4	Allcott (2011)	Electricity consumption	Report	40,000	Yes	No	No	No	Professor X or a company
5	Hanna, Duflo and Greenstone (2016)	Health status	Improved cooking stove	1,600	No	No	Yes	Yes	Professor X or an NPO
6	Hosono and Aoyagi (2018)	Sorting waste	Opportunity to win a laundry detergent	500	No	Yes	No	Yes	Professor X or an IDA

*Notes:* This table summarizes the six experiments examined in the present study. Column 9 presents the implementer of the program in each description.

Table 2: Summary statistics of the first online survey

	Mean (1)	SD (2)
Female	0.480	0.500
Age	46.673	14.064
Married	0.609	0.488
Living with children	0.379	0.485
Household income (10 thousand JPY)	535.289	249.164
Full-time employee	0.249	0.432
Temporary/contract employee	0.052	0.222
Self-employed	0.056	0.229
Part-time employee	0.124	0.330
Housewife/househusband	0.181	0.385
Unemployed/retired	0.103	0.305
Lives in Tokyo	0.125	0.331
Lives in Osaka	0.072	0.258

*Notes:* This table reports the means and standard deviations from the first survey. The number of observations is 2,107, except for Household income (the number of observations is 1,645).



Table 3: Comparisons of ethical concerns in the six studies (coefficients)

	Ordered logit		OLS	
	(1)	(2)	(3)	(4)
Fryer et al. (2015)	-0.418*** (0.090)	-0.420*** (0.090)	-0.219*** (0.048)	-0.219*** (0.047)
Thornton (2008)	0.046 (0.088)	0.036 (0.087)	0.037 (0.047)	0.032 (0.046)
Landry et al. (2006)	0.663*** (0.087)	0.665*** (0.087)	0.367*** (0.047)	0.366*** (0.046)
Allcott (2011)	-0.276*** (0.097)	-0.283*** (0.096)	-0.129** (0.051)	-0.136*** (0.051)
Hanna et al. (2016)	0.072 (0.110)	0.071 (0.110)	0.044 (0.059)	0.043 (0.058)
Order (1-6)	-0.069*** (0.013)	-0.070*** (0.013)	-0.034*** (0.007)	-0.034*** (0.007)
Female		0.304*** (0.074)		0.163*** (0.039)
Age		0.011*** (0.003)		0.006*** (0.001)
Constant			3.197*** (0.046)	2.830*** (0.082)
Control implementers	Yes	Yes	Yes	Yes
Other control variables	No	Yes	No	Yes
Number of Observations	6321	6321	6321	6321
Pseudo R-squared / R-squared	0.013	0.019	0.037	0.054

*Notes:* This table reports the estimates from the regression analyses in which the dependent variable is the response to the question “Do you recognize any ethical issues in this study?” on a five-point scale (1–5), as shown in Figure 1. Five studies are compared to Hosono and Aoyagi (2018). The coefficients are reported. Standard errors, clustered at the respondent level, are in parentheses. Columns 2 and 4 include other variables in Table 2 as well as Time spent on the survey. \*\*\*, \*\*, and \* indicate significance at the 1%, 5%, and 10% levels, respectively.

Table 4: Results of the randomized survey experiment (Fryer et al., 2015)

	Ordered logit		OLS	
	(1)	(2)	(3)	(4)
T1: Deleting the informed consent statement	0.312*** (0.114) [0.007]	0.312*** (0.114) [0.007]	0.173*** (0.059) [0.004]	0.172*** (0.059) [0.004]
T2: Samples are selected rather than self-selection	0.273** (0.115) [0.016]	0.271** (0.115) [0.016]	0.156*** (0.061) [0.009]	0.154** (0.061) [0.010]
T3: Deleting holiday parties statement	0.005 (0.113) [0.963]	0.006 (0.113) [0.963]	0.008 (0.058) [0.887]	0.006 (0.058) [0.907]
Order (1/2)		-0.193** (0.080)		-0.078* (0.041)
Constant			2.840*** (0.043)	2.957*** (0.073)
Multiple-Hypothesis Testing	0.018	0.018	0.010	0.011
Number of Observations	2146	2146	2146	2146
Pseudo R-squared / R-squared	0.002	0.003	0.007	0.009

*Notes:* This table reports the estimates from regression analyses in which the dependent variable is the response to the question “Do you recognize any ethical issues in this study?” on a five-point scale (1–5). The impact of changing the description on Fryer et al. (2015) to three treatment descriptions is evaluated. The coefficients are reported. Standard errors are in parentheses in columns 1 and 2. Robust standard errors are in parentheses in columns 3 and 4. The randomization- $t$   $p$ -values are in brackets. Inference in each column is based on a randomization inference procedure of Young (2019). \*\*\*, \*\*, and \* indicate significance at the 1%, 5%, and 10% levels, respectively. The row of Multiple-Hypothesis Testing reports the randomization- $t$   $p$ -values for the multiple-hypothesis testing test based on a randomization inference procedure of Young (2019), which applies the procedure of Westfall and Young (1993). It tests the null hypothesis that all treatment effects in each equation (each column) are zero.

Table 5: Results of the randomized survey experiment (Landry et al., 2006)

	Ordered logit		OLS	
	(1)	(2)	(3)	(4)
T1: Before-after study without control	-0.081 (0.112) [0.478]	-0.075 (0.112) [0.515]	-0.058 (0.061) [0.350]	-0.054 (0.061) [0.386]
T2: Treatment is a message rather than a raffle	-0.181* (0.111) [0.097]	-0.176 (0.111) [0.109]	-0.087 (0.059) [0.132]	-0.085 (0.059) [0.142]
T3: Promoting waste sorting rather than donations	-0.396*** (0.111) [0.000]	-0.400*** (0.111) [0.000]	-0.208*** (0.060) [0.000]	-0.211*** (0.059) [0.000]
Order (1/2)		-0.143* (0.079)		-0.087** (0.043)
Constant			3.452*** (0.042)	3.584*** (0.074)
Multiple-Hypothesis Testing	0.001	0.001	0.001	0.001
Number of Observations	2146	2146	2146	2146
Pseudo R-squared / R-squared	0.002	0.003	0.006	0.008

*Notes:* This table reports the estimates from regression analyses in which the dependent variable is the response to the question “Do you recognize any ethical issues in this study?” on a five-point scale (1–5), as shown in Figure 3. The impact of changing the description on Landry et al. (2006) to three treatment descriptions is evaluated. The coefficients are reported. Standard errors are in parentheses in columns 1 and 2. Robust standard errors are in parentheses in columns 3 and 4. The randomization- $t$   $p$ -values are in brackets. Inference in each column is based on a randomization inference procedure of Young (2019). \*\*\*, \*\*, and \* indicate significance at the 1%, 5%, and 10% levels, respectively. The row of Multiple-Hypothesis Testing reports the randomization- $t$   $p$ -values for the multiple-hypothesis testing test based on a randomization inference procedure of Young (2019), which applies the procedure of Westfall and Young (1993). It tests the null hypothesis that all treatment effects in each equation (each column) are zero.

Table 6: Results of subsample analyses

	Female			Male		
	(1)	(2)	(3)	(4)	(5)	(6)
	Ologit	Ologit	OLS	Ologit	Ologit	OLS
T1: Deleting the informed consent statement	0.645*** (0.165) [0.000]	0.646*** (0.166) [0.000]	0.329*** (0.081) [0.000]	0.031 (0.157) [0.840]	0.032 (0.157) [0.839]	0.022 (0.085) [0.796]
T2: Samples are selected rather than self-selection	0.448*** (0.163) [0.006]	0.433*** (0.162) [0.007]	0.230*** (0.078) [0.004]	0.100 (0.163) [0.535]	0.101 (0.163) [0.531]	0.069 (0.093) [0.452]
T3: Deleting holiday parties statement	-0.030 (0.165) [0.859]	-0.032 (0.164) [0.849]	-0.007 (0.080) [0.927]	0.030 (0.155) [0.847]	0.031 (0.155) [0.844]	0.015 (0.084) [0.858]
Order (1/2)		-0.358*** (0.116)	-0.163*** (0.056)		-0.040 (0.111)	0.005 (0.061)
Constant			3.056*** (0.097)			2.861*** (0.108)
Multiple-Hypothesis Testing	0.001	0.001	0.001			
Number of Observations	1052	1052	1052	1094	1094	1094
Pseudo R-squared / R-squared	0.009	0.012	0.033	0.000	0.000	0.001

	Female			Male		
	(1)	(2)	(3)	(4)	(5)	(6)
	Ologit	Ologit	OLS	Ologit	Ologit	OLS
T1: Before-after study without control	-0.064 (0.163) [0.703]	-0.065 (0.163) [0.697]	-0.033 (0.085) [0.705]	-0.105 (0.154) [0.492]	-0.089 (0.154) [0.559]	-0.074 (0.087) [0.394]
T2: Treatment is a message rather than a raffle	-0.367** (0.159) [0.020]	-0.365** (0.159) [0.021]	-0.180** (0.080) [0.023]	-0.035 (0.156) [0.825]	-0.026 (0.156) [0.870]	0.002 (0.086) [0.979]
T3: Promoting waste sorting rather than donations	-0.568*** (0.161) [0.001]	-0.576*** (0.161) [0.000]	-0.295*** (0.081) [0.000]	-0.262* (0.155) [0.090]	-0.262* (0.155) [0.092]	-0.132 (0.087) [0.128]
Order (1/2)		-0.123 (0.114)	-0.091 (0.057)		-0.154 (0.111)	-0.081 (0.062)
Constant			3.693*** (0.103)			3.478*** (0.107)
Multiple-Hypothesis Testing	0.003	0.003	0.002			
Number of Observations	1052	1052	1052	1094	1094	1094
Pseudo R-squared / R-squared	0.006	0.006	0.018	0.001	0.002	0.005

*Notes:* This table reports the estimates from subsample analyses of Tables 4 and 5. The coefficients are reported. Standard errors are in parentheses in columns 1, 2, 4 and 5. Robust standard errors are in parentheses in columns 3 and 6. The randomization- $t$   $p$ -values are in brackets. Inference in each column is based on a randomization inference procedure of Young (2019). \*\*\*, \*\*, and \* indicate significance at the 1%, 5%, and 10% levels, respectively. The row of Multiple-Hypothesis Testing reports the randomization- $t$   $p$ -values for the multiple-hypothesis testing test based on a randomization inference procedure of Young (2019), which applies the procedure of Westfall and Young (1993). For example, column 1 reports the result under the null hypothesis that all treatment effects within the two regressions of the same model for women (column 1) and men (column 4) are zero, adjusting for multiple-hypothesis testing.