

# Behavioural spillovers unpacked: estimating the side effects of social norm nudges

Julien Picard and Sanchayan Banerjee

October 2023

Centre for Climate Change Economics  
and Policy Working Paper No. 427  
ISSN 2515-5709 (Online)

Grantham Research Institute on  
Climate Change and the Environment  
Working Paper No. 402  
ISSN 2515-5717 (Online)

**The Centre for Climate Change Economics and Policy (CCCEP)** was established by the University of Leeds and the London School of Economics and Political Science in 2008 to advance public and private action on climate change through innovative, rigorous research. The Centre is funded by the UK Economic and Social Research Council. Its third phase started in October 2018 with seven projects:

1. Low-carbon, climate-resilient cities
2. Sustainable infrastructure finance
3. Low-carbon industrial strategies in challenging contexts
4. Integrating climate and development policies for 'climate compatible development'
5. Competitiveness in the low-carbon economy
6. Incentives for behaviour change
7. Climate information for adaptation

More information about CCCEP is available at [www.cccep.ac.uk](http://www.cccep.ac.uk)

**The Grantham Research Institute on Climate Change and the Environment** was established by the London School of Economics and Political Science in 2008 to bring together international expertise on economics, finance, geography, the environment, international development and political economy to create a world-leading centre for policy-relevant research and training. The Institute is funded by the Grantham Foundation for the Protection of the Environment and a number of other sources. It has 13 broad research areas:

1. Biodiversity
2. Climate change adaptation and resilience
3. Climate change governance, legislation and litigation
4. Climate, health and environment
5. Environmental behaviour
6. Environmental economic theory
7. Environmental policy evaluation
8. International climate politics
9. Science and impacts of climate change
10. Sustainable public and private finance
11. Sustainable natural resources
12. Transition to zero emissions growth
13. UK national and local climate policies

More information about the Grantham Research Institute is available at: [www.lse.ac.uk/GranthamInstitute](http://www.lse.ac.uk/GranthamInstitute)

**Suggested citation:**

Picard J and Banerjee S (2023) *Behavioural spillovers unpacked: estimating the side effects of social norm nudges*. Centre for Climate Change Economics and Policy Working Paper 427/Grantham Research Institute on Climate Change and the Environment Working Paper 402. London: London School of Economics and Political Science

# Behavioural Spillovers Unpacked:

## Estimating the Side Effects of Social Norm Nudges<sup>†</sup>

Julien Picard<sup>a\*</sup>, Sanchayan Banerjee<sup>b</sup>

(Last version available [here](#))

### Abstract

Fighting the climate crisis requires changing many aspects of our consumption habits. Previous studies show that a first climate-friendly action can lead to another. Does deciding not to eat meat increase our willingness to do more for the environment? Can encouraging vegetarianism alter this willingness? Using an online randomised control trial, we study the side effects of a social norm nudge promoting vegetarianism on environmental donations. We develop an experimental design to estimate these side effects and a utility maximisation framework to understand their mechanisms. Using an instrumental variable, we find that choosing not to eat meat increases donations to pro-environmental charities. We use machine learning to find that the social norm nudge crowds out donations from the population segment prone to choosing vegetarian food after seeing the nudge. However, the nudge led another group to make less carbon-intensive food choices without affecting their donations. Our results suggest that whilst social norm nudges are effective on specific population segments, they can also reduce the willingness of some groups to do more. *JEL* codes: C36, C93, D91, Z18.

**Keywords:** social norm; meat; climate change; behavioural spillovers; side effects

**CRedit author statement (details [here](#)):**

JP: conceptualisation, data curation, formal analysis, investigation, funding acquisition, software, methodology, project administration, visualisation, writing.

SB: funding acquisition.

<sup>†</sup>We are particularly grateful to Eugenie Dugoua, Marion Dumas, Elisabeth Gsottbauer and Gregor Singer for their invaluable advice. We also thank Susana Mourato, Kelsey Jack, Christopher Roth, Jeanne Hagenbach, Mark Andor, Meryem Yankol-Schalck, Sarah Gharbi and many conference participants for their insightful comments on the paper. We also thank the Royal Geographic Society for its financial support. Finally, we acknowledge support in publishing this paper from the Grantham Foundation for the Protection of the Environment and the Economic and Social Research Council through the Centre for Climate Change Economics and Policy

<sup>a</sup>Department of Geography and Environment, London School of Economics, London.

<sup>b</sup>Department of Environmental Economics, Vrije Universiteit Amsterdam.

\*Corresponding author: Julien, Picard; [j.r.picard@lse.ac.uk](mailto:j.r.picard@lse.ac.uk).

# I Introduction

Climate change is undoubtedly one of the most critical challenges of the 21<sup>st</sup> century, with potentially devastating economic consequences if urgent and ambitious actions are not taken. One way to curb emissions is to make significant lifestyle changes, especially in developed countries (Shukla et al., 2022). In this regard, research suggests that an initial pro-environmental act can influence our propensity to engage in subsequent ones. However, this “behavioural spillover effect,” as termed by Thøgersen (1999), can manifest in multiple ways. For example, Comin and Rode (2023) provide causal evidence that people who install solar panels are more likely to vote for green parties, whilst Mazar and Zhong (2010) note a decline in altruism among those who purchase green products. Therefore, policymakers should encourage changes in our everyday actions that not only significantly reduce carbon emissions but also inspire people to do even more for the environment.

Cutting on meat yields large reductions in greenhouse gas emissions (Green et al., 2015; Riahi et al., 2022). However, does reducing meat also increase our inclination to do more? And would encouraging meat reduction with social norm nudges alter this inclination? Social norm nudges are simple messages giving information on what others are doing, approve or disapprove (Bicchieri, 2016). These messages are effective in shifting behaviours<sup>1</sup> and are quite popular in the environmental domain. They have been used to foster recycling,<sup>2</sup> promote sustainable diets,<sup>3</sup> improve water and electricity consumption,<sup>4</sup> and even foster towel reuse in

---

<sup>1</sup>See Rhodes, Shulman, and McClaran (2020) and Melnyk, van Herpen, Trijp, et al. (2010) for meta-analyses on the effectiveness of social norm messaging in general. For meta-analyses and reviews of the effectiveness of social norm messaging applied to the environmental domain, see Abrahamse and Steg (2013); Andor and Fels (2018); Cialdini and Jacobson (2021); Farrow, Grolleau, and Ibanez (2017).

<sup>2</sup>See for instance Andersson and von Borgstede (2010); Bratt (1999); Fornara, Carrus, Passafaro, and Bonnes (2011); Nigbur, Lyons, and Uzzell (2010).

<sup>3</sup>See for instance Richter, Thøgersen, and Klöckner (2018); Salmivaara and Lankoski (2019); Sparkman and Walton (2017); Sparkman, Weitz, Robinson, Malhotra, and Walton (2020); Stea and Pickering (2019); Testa, Russo, Cornwell, McDonald, and Reich (2018); Wenzig and Gruchmann (2018).

<sup>4</sup>See for instance Allcott (2011); Carrico and Riemer (2011); Costa and Kahn (2013); Ferraro, Miranda, and Price (2011); Handgraaf, De Jeude, and Appelt (2013); Lapinski, Rimal, DeVries, and Lee (2007); Nolan, Schultz, Cialdini, Goldstein, and Giskevicius (2008).

hotels.<sup>5</sup> Yet, little is known about whether these messages trigger any side effects on other non-targeted pro-environmental decisions.

In this paper, we study the side effects of a social norm message promoting vegetarianism in an online experiment with 2,775 English respondents. Respondents choose their preferred meal on a restaurant menu. At the end of the experiment, respondents can donate to a pro-environmental charity of their choice. We use this task to proxy their willingness to make more effort for the environment.

To understand the mechanisms underpinning the side effects of social norm nudges, we model behavioural spillovers and the effect of behavioural policies in a utility maximisation framework. To our knowledge, our model is the first to formalise why acting pro-environmentally can increase or decrease our willingness to do extra climate-friendly deeds and the mechanisms leading policy interventions to alter this willingness.<sup>6</sup> We assume that the trajectory of behavioural spillovers hinges on the motivation behind the initial pro-environmental act. On the one hand, acting pro-environmentally as an *end to itself* (e.g. because it is "something we care about") raises the cost of reverting to self-serving behaviours (positive behavioural spillovers). For instance, quitting meat due to conviction may encourage us to make other eco-friendly habits to stay true to our beliefs. On the other hand, acting pro-environmentally as *a means to an end* reduces the cost of reverting (negative behavioural spillovers). For instance, quitting meat out of social pressure can license us to do less because the initial action was not self-driven.

Behavioural policies that foster pro-environmental acts target intrinsic or extrinsic motives. By doing so, they also affect non-targeted green behaviours through two channels. The first is through their effect on the initial pro-environmental action, subsequently triggering behavioural spillovers. The second is direct through their effect on motivations. We term this latter influence a crowding-in/out effect. The sign of crowding-in/out effects indicates the

---

<sup>5</sup>See for instance [Goldstein, Cialdini, and Griskevicius \(2008\)](#); [Reese, Loew, and Steffgen \(2014\)](#); [Schultz, Khazian, and Zaleski \(2008\)](#).

<sup>6</sup>For an extended version of the model, see [Picard \(2021\)](#).

causal mechanisms through which policies operate. Crowding-in effects suggest that policies have induced individuals to act pro-environmentally as an end to itself. Crowding-out effects imply that policies have led individuals to act as a means to an end.

Estimating behavioural spillovers and policies' crowding-in/out effects is, therefore, crucial when studying the side effects of policies. Nevertheless, causally estimating behavioural spillovers is difficult. We embed an instrumental variable in our experimental design to causally estimate the effect of choosing a vegetarian meal on charitable giving. Namely, beyond merely allocating participants into control (no message) and treatment groups (receiving the social norm message), we randomly show them different selections of items varying in their meat content. Respondents have to opt out to access the full menu. Allocation to these selections alters the likelihood of choosing a vegetarian dish without directly affecting donations. We show that our instrumental variable strategy allows us to estimate the main equation of our model, capturing the relationship between the main effect of policies, behavioural spillovers, and the policy's crowding-in/out effect.

A second complication that arises when looking at the side effect of social norm messages is that they can be perceived differently depending on people's willingness to follow the norm. Therefore, the main effects of the social norm on food choices and their side effects on charitable giving may be heterogeneous. To explore this heterogeneity, we use data from another treatment arm (n=2782) where respondents reveal their willingness to transition to vegetarian diets after reading the social norm message. We use this extra survey data to train a gradient tree boosting classifier to predict this willingness based on respondents' social-demographic characteristics, attitudinal information and self-reported beliefs (Friedman, 2001). This method provides distinct advantages over approaches like mediation analysis or machine learning techniques, as detailed by Künzel, Sekhon, Bickel, and Yu (2019) or Athey and Imbens (2015). First, it does not directly rely on direct measurements, sidestepping the challenges of pre-treatment questions that may inadvertently hint at study objectives. Second, it enables the examination of non-linear associations without unnecessarily expanding the set of regressors.

Last, by defining profiles *ex-ante*, our method makes the source of heterogeneity explicit, i.e., people's willingness to conform to the norm.

Our results indicate that the social norm nudge is effective. Exposure to the message increases the likelihood of choosing a vegetarian item on average. Respondents predicted to be trying to follow the norm drive this effect. However, this does not translate into a significant decrease in the carbon footprint of food choices for this subsample. Conversely, respondents predicted to be hesitant about following the norm are not more likely to choose vegetarian food after seeing the nudge but emit significantly less with their choices. We do not observe an effect of the social norm nudge on respondents predicted to be unwilling to conform and on respondents predicted to be already conforming. Our experimental results provide insights into the specific social-demographic profiles that policymakers should target to maximise the effect of social norm messaging. In line with [Bryan, Tipton, and Yeager \(2021\)](#)'s recommendation to systematically address heterogeneity when evaluating behavioural policies, our results confirm the importance of doing so.

We also find evidence of a positive behavioural spillover effect on average. Namely, respondents choosing vegetarian food are more likely to give to pro-environmental charities. This indicates that choosing to eat less meat can encourage people to do more for the environment. However, the social norm message crowds out donations of those predicted to be trying to conform. This crowding-out effect dominates the positive behavioural spillover effect. Our model suggests that social norm messaging pushes this group to act pro-environmentally out of extrinsic motivation (e.g., to comply with social pressure). This, in turn, reduces their engagement in the donation task.

We contribute to a burgeoning literature studying the side effects of policies. [Maki et al. \(2019\)](#) and [Geiger, Brick, Nalborczyk, Bosshard, and Jostmann \(2021\)](#) identify only weak evidence for such side effects in their meta-analyses. Methodological discrepancies could explain

this scarcity of compelling evidence by making studies hard to compare with one another.<sup>7</sup> The instrumental variable strategy we use in this paper aligns with the recommendations made by Bonev (2023) for estimating behavioural spillover effects. To our knowledge, only Alacevich, Bonev, and Söderberg (2021), Comin and Rode (2023), and Alt, Bruns, and DellaValle (2023) have used an instrumental variable approach to estimating behavioural spillovers. Our paper is, however, the first to precisely assess the potential trade-offs between promoting pro-environmental behaviours and crowding out others. Finally, using an empirical strategy grounded in theory, our approach allows us to infer the cognitive mechanisms of policies from the sign of their crowding-in/out effects.

The remaining of this article is articulated as follows. We present our theoretical model in Section II. In Section III, we present our empirical framework. In Section IV, we share the experimental results from a social norm nudge promoting vegetarianism. Section V concludes.

## II Modelling Behavioural Spillovers

Positive behavioural spillovers are often explained using cognitive dissonance theory (Festinger, 1962). Cognitive dissonance theory posits that people prefer staying consistent across their choices. This is particularly true when pro-environmental deeds signal an altruistic identity, which raises the mental discomfort of behaving at odds with this identity. Conversely, moral licensing theory explains negative behavioural spillovers (Merritt, Effron, & Monin, 2010). In this case, a first pro-environmental deed allows decision-makers to subsequently act selfishly.

---

<sup>7</sup>Some studies compare respondents exposed to a policy with those allocated to a control group (Carrico, Raimi, Truelove, & Eby, 2018; Liu, Kua, & Lu, 2021; Van Rookhuijzen, De Vet, & Adriaanse, 2021; Wolstenholme, Poortinga, & Whitmarsh, 2020). This method does not distinguish policies' crowding-in/out effect from behavioural spillovers. Other studies randomly offer participants the targeted behaviours to estimate subsequent spillover effects (Alt & Gallier, 2022; Clot, Della Giusta, & Jewell, 2022; Margetts & Kashima, 2017). This design supposes that choosing (not) to do the targeted behaviour is the same as (not) being proposed to do it. This assumption is, however, debatable.



In our model, the motivation underpinning the first pro-environmental deed determines the direction of behavioural spillovers. This allows us to reconcile the predictions made by these two theories in a utility maximisation framework. By modelling behavioural policies as altering these motivations, our theory captures how such policies interact with behavioural spillovers.

**Motivations:** We consider two motives for acting pro-environmentally. Individuals' decisions are intrinsically motivated or extrinsically motivated. Intrinsic motivations imply "identity-enhancing" decisions:<sup>8</sup> we act pro-environmentally because it is *who we are*. For instance, imagine the fictional character of Anne, who stops eating meat because she cares about the environment. She will signal a pro-environmental identity when eating with her colleagues. This may induce her to engage in other pro-environmental behaviours to avoid appearing inconsistent.<sup>9</sup> Experimental evidence shows that consistent behaviours are more frequent when reminding people of their past actions (e.g., [Gneezy, Imas, Brown, Nelson, and Norton 2012](#); [Lacasse 2016](#); [Van der Werff, Steg, and Keizer 2013, 2014](#)), after pledging (e.g., [Banerjee, Galizzi, John, and Mourato 2022](#); [Lokhorst, Werner, Staats, van Dijk, and Gale 2013](#)), or when labelling people as pro-environmental (e.g., [Baca-Motes, Brown, Gneezy, Keenan, and Nelson 2013](#); [Lacasse 2016](#)).

On the other hand, extrinsic motivations imply decisions done as *a means to an end*. Extrinsic motivations range from seeking material rewards (e.g., a tax rebate) to more intangible rewards (e.g., not feeling excluded). In such instances, moral licensing is more likely to occur. To illustrate this, imagine now the fictional character of Bob, a colleague of Anne, who starts choosing the vegetarian option at the cafeteria because he feels pressured to do so.

---

<sup>8</sup>Following [Akerlof and Kranton \(2000\)](#), we define *identity* as a sense of self. To enhance this sense of self, one has to engage in certain behaviours (e.g., good Samaritans help people, Harley Davidson bikers prefer beer over hot milk, and environmentalists sort their waste).

<sup>9</sup>As observed early on by Adam Smith, one does not need to be scrutinised by others to fear appearing inconsistent: "When I endeavour to examine my own conduct, [...] I divide myself, as it were, into two persons: and that I, the examiner and judge, represent a different character from that other I, the person whose conduct is examined and judged of" ([Smith, 1853](#)).

This pro-environmental deed does not enhance his pro-environmental identity: Bob acts pro-environmentally to avoid disapproving looks. Once the vegetarian choice is made, Bob can feel relieved, think he has done his bit and indulge in other self-serving behaviours. In the literature, financial incentives have been found to induce people to slacken after exercising a certain level of effort (e.g., [Dolan and Galizzi 2014](#); [Hartmann, Marcos, and Barrutia 2023](#); [Steinhorst and Matthies 2016](#); [Xu, Zhang, and Ling 2018](#)). Similarly, social pressure can also induce moral licensing (e.g., [Kristofferson, White, and Peloza 2014](#); [Tiefenbeck, Staake, Roth, and Sachs 2013](#)).

**Context:** The context in which people make decisions influences the relative importance of each motive. Educative pieces of information can induce intrinsically motivated pro-environmental deeds, for instance. Conversely, situations where rewards are contingent on acting pro-environmentally can induce extrinsically motivated actions. Thus, how our past choices influence our current ones depends on the context in which we made these choices. In this model, we assume that individuals perfectly remember the context in which they make their decisions.

**Decision utility:** Let us define the function describing individuals' decision processes. We consider a simple two-period model. We assume that decision-makers only consider the current period when making a choice. Their decision utility at period 2 is of the form:

$$U_2 \equiv u(x_2|I_1, E_1, \eta_2, \epsilon_2) - c_2 \cdot x_2 \quad (\text{II.1})$$

$x_2$  is the amount of effort spent acting pro-environmentally in period 2. In our experiment, this corresponds to donations to pro-environmental charity, whilst  $x_1$  corresponds to the efforts made for the environment when choosing food.  $\eta_2$  denotes the propensity of individuals to act out of intrinsic motivation.  $\epsilon_2$  denotes the propensity of individuals to act out of extrinsic motivation. Decision utility (II.1) is increasing and concave in  $x_2$ ,  $\eta_2$  and  $\epsilon_2$ . The marginal

utility of doing pro-environmental deeds  $x_2$  also increases in  $\eta_2$  and  $\epsilon_2$ .  $c_2$  is the cost of exerting one unit of pro-environmental effort. It captures the difficulty of acting pro-environmentally (e.g., the number of steps before making an online donation to a charity). The context in which decisions are made alters parameters  $c_2$ ,  $\eta_2$  and  $\epsilon_2$ .

Functions  $I_1$  and  $E_1$  capture the influence of choices of period 1 on choices of period 2. They represent respectively the extent to which individuals remember being intrinsically or extrinsically motivated, such that  $I_1 : \{x_1, \eta_1\} \mapsto \mathbb{R}^+$  and  $E_1 : \{x_1, \epsilon_1\} \mapsto \mathbb{R}^+$ . Both functions are assumed to be increasing and concave in their arguments with positive cross-derivatives. The latter assumption implies that the higher the efforts spent on acting pro-environmentally, and the higher  $\eta_1$  ( $\epsilon_1$ ), the more individuals remember their decisions to be intrinsically (extrinsically) motivated. The decision utility of period one is defined similarly, with  $I_0$  and  $E_0$  given. In what follows, we focus on the effect of pro-environmental decisions of period one and the context in which they are made on period two pro-environmental decisions. We ignore the effect of the context of period two (i.e.,  $c_2$ ,  $\eta_2$  and  $\epsilon_2$ ) on period two choices. Let us define two competing mechanisms influencing the marginal utility of doing  $x_2$ : *consistency* and *moral licensing*.

**Definition 1.** *Consistency describes an increase in pro-environmental effort following a first pro-environmental deed. It occurs when remembering that one was intrinsically motivated in period 1 increases the utility of doing  $x_2$  ( $\partial_{I_1} U_2 > 0$  and  $\partial_{x_2 I_1} U_2 > 0$ ).*

**Definition 2.** *Moral licensing describes a decrease in pro-environmental effort following a first pro-environmental deed. It occurs when remembering that one was extrinsically motivated in period 1 reduces the utility of doing  $x_2$  ( $\partial_{E_1} U_2 < 0$  and  $\partial_{x_2 E_1} U_2 < 0$ ).*

**Main effect of behavioural policies** We assume a social planner seeking to increase individuals' efforts to act pro-environmentally. To do this, she designs a policy altering the decision-making context in period one. This policy can affect individuals' motivations by altering parameters  $\eta_1$  and  $\epsilon_1$  or reduce the cost of doing the targeted pro-environmental behaviour

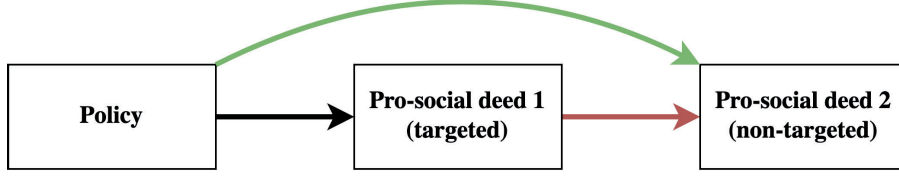


Figure 1: Side effects a policy

Note: causal mechanisms of the effects of a policy on non-targeted decisions. The red arrow represents behavioural spillovers, whilst the green arrow represents crowding-out/in effects.

$c_1$ . We refer to policies altering parameters  $\eta_1$  and  $\epsilon_1$  as *preference nudges* and policies playing on  $c_2$  as *choice architecture nudges*.<sup>10</sup>

**Lemma 1.** *In period one, behavioural policies increase pro-environmental efforts when increasing  $\theta_1 \in \{\eta_1, \epsilon_1\}$  or when reducing  $c_1$ .*

$$\frac{\partial x_1^*}{\partial \theta_1} = \frac{\partial_{x_1 \theta_1} U_1}{-\partial_{x_1 x_1} U_1} > 0 \quad \frac{\partial x_1^*}{\partial c_1} = \frac{1}{\partial_{x_1 x_1} U_1} < 0 \quad (\text{II.2})$$

See the proof in Appendix A.  $x_1^*$  is the amount of effort allocated to pro-environmental activities maximising individuals' decision utility.

**Side effects of policies:** *Preference nudges* affect the non-targeted decision of period 2 through two channels as described by equation (II.3). The first is through the effect of the behaviour targeted by the policy on the non-targeted behaviour, which we refer to as a "behavioural spillover". The red arrow represents this "behavioural spillover" in Figure 1. The second channel is through the effect of the policy on individuals' motivations. The green arrow represents this "crowding-in/out" effect" in Figure 1.

**Lemma 2.** *Preference nudges  $\theta_1 \in \{\eta_1, \epsilon_1\}$  alter decisions of period two through two channels captured by the following equation:*

$$\underbrace{\frac{dx_2^*}{d\theta_1}}_{\text{Side effect}} = \frac{1}{-\partial_{x_2 x_2} U_2} \left( \underbrace{\partial_{x_1 x_2} U_2}_{\text{Behavioural spillover}} \times \underbrace{\frac{\partial x_1^*}{\partial \theta_1}}_{\text{Main effect}} + \underbrace{\partial_{x_2 \theta_1} U_2}_{\text{Crowding-in/out effect}} \right) \quad (\text{II.3})$$

<sup>10</sup>Here, we follow Löfgren and Nordblom (2020)'s typology.

Choice architecture nudges only affect decisions of period 2 through a behavioural spillover effect.

**Lemma 3.** *Choice architecture nudges only alter decisions of period 2 through one channel captured by the following equation:*

$$\underbrace{\frac{dx_2^*}{dc_1}}_{\text{Side effect}} = \underbrace{\frac{\partial_{x_1x_2}U_2}{-\partial_{x_2x_2}U_2}}_{\text{Behavioural spillover}} \times \underbrace{\frac{\partial x_1^*}{\partial c_1}}_{\text{Main effect}} \quad (\text{II.4})$$

See Appendix A for proofs of lemma 2 and 3. The sign of the behavioural spillover effect indicates whether individuals are more intrinsically motivated (and therefore more consistent) than they are extrinsically motivated (prone to engage in moral accounting):

$$\partial_{x_1x_2}U_2 = \underbrace{\partial_{x_2I_1}U_2 \cdot \partial_{x_1}I_1}_{\text{Consistence } (>0)} + \underbrace{\partial_{x_2E_1}U_2 \cdot \partial_{x_1}E_1}_{\text{Moral licensing } (<0)} \quad (\text{II.5})$$

It, therefore, captures the degree of complementarity between the targeted and the non-targeted pro-environmental decisions. Positive (negative) behavioural spillover effects indicate that  $x_1$  and  $x_2$  are complements (substitutes): doing  $x_1$  increases (decreases) the marginal utility of doing  $x_2$ . The crowding-in/out effect captures the effect of the policy on period two choices through its impact on individuals' motivations:

$$\partial_{x_2\eta_1}U_2 = \underbrace{\partial_{x_2I_1}U_2 \cdot \partial_{\eta_1}I_1}_{\text{Crowding-in}}, \quad \partial_{x_2\varepsilon_1}U_2 = \underbrace{\partial_{x_2E_1}U_2 \cdot \partial_{\varepsilon_1}E_1}_{\text{Crowding-out}} \quad (\text{II.6})$$

When playing on intrinsic (extrinsic) motivations, policies increase (decrease) the marginal utility of doing  $x_2$ .

**Proposition 1.** *Behavioural spillovers and crowding-in/out effects can have opposite signs.*

The sign of the behavioural spillover effect does not depend on the policy. Conversely, the crowding-in/out effect is only positive when policies affect individuals' intrinsic motivations. Therefore, favouring policies that increase people's intrinsic motivations is crucial to strengthening people's willingness to do extra pro-environmental actions.

**Proposition 2.** *Playing on intrinsic motivations maximises the side effects of policies.*

In the next section, we develop an experimental design which allows us to estimate these two effects.

### III Empirical Strategy

Behavioural spillovers capture whether people perceive two actions as complements or substitutes. Estimating this effect is crucial when policymakers seek to promote behaviours that serve as “entry points” to other pro-environmental actions. Additionally, crowding-in/out effects inform us of whether policies affect people’s intrinsic or extrinsic motivations. However, estimating behavioural spillovers and crowding-in/out effects is not trivial. Two complications arise. First, causally estimating behavioural spillovers is difficult: unobserved variables can be correlated with people’s participation in the targeted and the non-targeted behaviours (e.g., values, beliefs). Second, the same policy can play on intrinsic motivations for one population segment and extrinsic motivations for another. Therefore, different crowding-in/out effects can emerge from one person to another. This section develops an empirical framework to address these two issues.

**Addressing omitted variable biases** First, we assume a population of  $N$  individuals indexed by  $i$  randomly exposed to a policy fostering a given pro-environmental deed. Denote by  $\mathbf{x}_1$  and  $\boldsymbol{\theta}_1$  the  $N \times 1$  vectors capturing respectively the effort spent by individuals on the targeted pro-environmental deed and their treatment status. The following linear models estimate respectively the main effect of the policy and its side effects on a non-targeted pro-environmental deed  $\mathbf{x}_2$ :

$$x_{1i} = \alpha^{ME} + \beta^{ME}\theta_{1i} + \varepsilon_i^{ME} \quad (\text{III.1})$$

$$x_{2i} = \alpha^{SE} + \beta^{SE}\theta_{1i} + \varepsilon_i^{SE} \quad (\text{III.2})$$

Here,  $\hat{\beta}^{ME}$  and  $\hat{\beta}^{SE}$  are, respectively, estimates of the main and the side effects of the policy. They are unbiased if the policy is randomised, i.e., if the error terms  $\varepsilon_i^{ME}$  and  $\varepsilon_i^{SE}$  are such that  $cov(\varepsilon^{ME}, \theta_1) = cov(\varepsilon^{SE}, \theta_1) = 0$ . As we showed in Section II, the side effect of policies can be decomposed into a behavioural spillover effect and a crowding-in/out effect. In what follows, we make the following assumption:

**Assumption 1.** *The magnitude and the sign of the behavioural spillover effect do not depend on the policy.*

Assumption 1 reflects the insights provided by the model. A naive approach to dissociate behavioural spillovers from policies' crowding-in/out effects consists of fitting the following linear model:

$$x_{2i} = \tilde{\alpha} + \tilde{\beta}^{BS} x_{1i} + \tilde{\beta}^C \theta_{1i} + \tilde{\varepsilon}_i \quad (\text{III.3})$$

Where  $\hat{\beta}^{BS}$  and  $\hat{\beta}^C$  are respectively estimates of behavioural spillovers and crowding-in/out effects. However, individuals acting pro-environmentally may be different from those who do not. When the experimenter can not observe these differences, this creates an omitted variable bias:  $cov(x_{1i}, \tilde{\varepsilon}_i) \neq 0$  implies that  $\hat{\beta}^{BS}$  and  $\hat{\beta}^C$  are biased. Our proposed solution is to instrument the pro-environmental decision of period one by any element in the choice architecture that alters the difficulty of undertaking the targeted decision without changing people's motivations. In an experimental design, this is equivalent to randomly allocating people to a pure *choice-architecture* nudge. Denote by  $\mathbf{c}_1$  the  $N \times 1$  vector capturing people's treatment status for this *choice architecture* nudge. A two-stage least square approach allows us to derive unbiased estimates of behavioural spillovers and crowding-in/out effects:

$$\begin{aligned} \text{Stage 1: } x_{1i} &= \alpha + \beta_1 c_{1i} + \beta_2 \theta_{1i} + \varepsilon_i \\ \text{Stage 2: } x_{2i} &= \alpha' + \beta^{BS} \hat{x}_{1i} + \beta^C \theta_{1i} + \varepsilon'_i \end{aligned} \quad (\text{III.4})$$

Where  $\hat{x}_{1i}$  are the predicted values for the first stage. For this approach to be valid, our instrumental variable should be relevant ( $cov(\mathbf{c}_1, \mathbf{x}_1) \neq 0$ ) and exogenous ( $cov(\mathbf{c}_1, \varepsilon') = 0$ ). We

obtain unbiased estimates of behavioural spillovers and crowding-in/out effects when this is the case. Furthermore, one can derive the following proposition:

**Proposition 3.** *Estimates of models (III.2) and (III.4) are such that:*

$$\underbrace{\hat{\beta}^{SE}}_{\text{Side effect}} = \underbrace{\hat{\beta}^{BS}}_{\text{Behavioural spillover}} \times \underbrace{\hat{\beta}^{ME}}_{\text{Main effect}} + \underbrace{\hat{\beta}^C}_{\text{Crowding in/out effect}} \quad (\text{III.5})$$

See Appendix A for the proof. Proposition 3 shows that our instrumental strategy allows us to interpret estimates of model (III.2) and (III.4) with equation (II.3) derived in Section II.

**Addressing heterogeneity** Different people may react differently to an intervention. We propose to explore this heterogeneity by defining different types expected to respond differently to the policy. We then use two separate data sets. An *inference* sample, where respondents are randomly allocated to a policy  $\theta_1$  and a choice architecture nudge  $c_1$ , and a *training* sample where new respondents reveal their types. We use the *training* data to train a classification algorithm to predict the types of respondents in the *inference* sample.

Put formally, let us index by  $j \in [1, \dots, N']$  the  $N'$  observations in the *training* sample where each observation's type  $y_j$  is known. Denote respectively by  $\mathbf{W}$  and  $\mathbf{W}'$  the  $N \times M$  and  $N' \times M$  matrices of covariates of the *inference* and the *training* samples. In three steps, we estimate the conditional average treatment effects of policy  $\theta_1$ . First, estimate the function  $y_i = f(W'_i)$  such that:

$$\hat{f} \in \arg \min_f L(y_i, f(W'_i)) \quad (\text{III.6})$$

Where  $L(\cdot)$  is a loss function. Then, predict the types of observations in the *inference sample*:

$$\hat{y}_i = \hat{f}(W_i) \quad (\text{III.7})$$

Finally, estimate treatment effects for each type.

The following section presents an application of this empirical framework to the case of a



social norm nudge promoting vegetarianism.

## IV Application

In the remainder of this article, we study the side effects of a social norm nudge promoting vegetarian diets in an online experiment delivered to 2,775 English respondents. In subsection [IV.A](#), we detail the design of this experiment. Then, subsection [IV.B](#) presents our empirical framework applied to our case of interest. Results are presented in subsection [IV.C](#). Finally, we explore heterogeneity in subsection [IV.D](#) using the method developed in [Section III](#).

### IV.A Experimental Design and Data Collection

We designed the survey experiment on Qualtrics and recruited respondents via Prolific. The experiment lasted approximately 10 minutes. The flat payment for completing it corresponds to Prolific’s standard payment rate, £5 per hour. Upon finishing the survey, respondents have a 1/100 chance to win a £20 voucher. In total, we recruited a sample of 5,557 English respondents divided between an *inference sample* (n=2,775) and a *training sample* (n=2,782).

We use the answers from the participants in the *training sample* to look at the heterogeneity in our treatment effects as part of an exploratory analysis (see subsection [IV.D](#)). Respondents in the *inference sample* participate in the main experiment, whose timeline is presented in [Figure 2](#).<sup>11</sup>

---

<sup>11</sup>The survey questionnaire can be found [here](#). We pre-registered the experimental design, power analysis, empirical strategy and instrumental variable strategy on Open Science Framework ([here](#)). The pre-analysis plan describes a broader project where three strands of research are investigated: 1) the effect of familiar food choices on one’s inclination to choose vegetarian food; 2) the effect of reflection on the effectiveness of social norm nudges (now published, [Banerjee and Picard 2023](#)); 3) the present study. When reporting our results, we correct for the pre-registered hypotheses. Deviations from the pre-analysis plan are documented and justified [here](#).

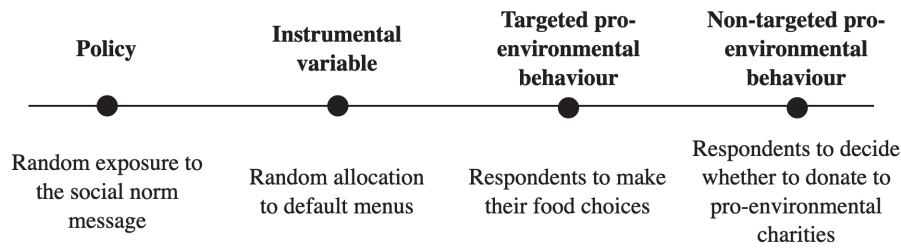


Figure 2: Timeline of the experiment

**Policy:** The policy of interest in this experiment is a social norm nudge. More precisely, we consider the following dynamic social norm message<sup>12</sup>:

*A study published in The Lancet Planetary Health found that the share of British people who stopped eating meat has increased by more than 50% from 2008 to 2019. More and more people are choosing plant-based dishes that are kinder to the planet and in turn, are becoming climate-friendly.*

Its formulation is similar to the one used by [Blondin, Attwood, Vennard, and Mayneris \(2022\)](#), who find that this message effectively increases vegetarian food choice intentions. In this experiment, we seek to determine whether choosing vegetarian food increases respondents' donations to a pro-environmental charity and whether such social norm messages crowd out or crowd in their contributions. As such, we randomly allocate respondents to a treatment group where this message is displayed (n=1391) or to a control group (n=1384) before making food choices.

**Instrumental variable:** As explained in Section III, omitted variables are likely to bias the estimation of behavioural spillovers and crowding-out/in effects. We tackle this issue by embedding an instrumental variable in the design: a default option nudge. Namely, when making their food choices, respondents first see a small menu containing a subset of food items presented as the chef's suggestion. Respondents are either randomly allocated to a default menu containing a higher share of meat-based items (n=1383, see Figure 6 in Appendix D) or a higher

<sup>12</sup>We construct it using the study of [Stewart, Piernas, Cook, and Jebb \(2021\)](#) analysing UK meat consumption trends using data from the National Diet and Nutrition Survey.

Table 1: Sample sizes of treatment groups

		Policy	
		Control	Treatment
Instrumental variable	Plant-intensive	690	693
	Meat-intensive	694	698

share of vegetarian items ( $n=1392$ , see Figure 5 in Appendix D). They can choose an item from this selection or opt out and access the main menu containing all the items. We expect that respondents are more likely to choose a vegetarian item when shown a selection containing mostly vegetarian items. Table 1 presents the sample size of each subgroup formed by the interaction between the random allocations to the policy and to default menus.

**Targeted pro-environmental behaviour:** We reproduce an online food order environment where participants choose a dish from a restaurant menu. Therefore, the targeted pro-environmental decision is whether participants choose vegetarian food.<sup>13</sup> We designed 24 different versions of the main menu, where we varied the items’ ordering and appearance. In all menus, we label food items with pictures of footprints ranging from green to red. An explanation indicates that green footprints mean “completely climate-friendly” and red footprints mean “not climate-friendly at all” (see Figure 4 in Appendix D). As such, all participants have the same baseline information regarding the environmental consequences of their choices. Table 21 in Appendix D displays the characteristics of the dishes shown in the menus.

<sup>13</sup>To mitigate any biases stemming from this choice being hypothetical, we ask two questions inspired by the literature on willingness-to-pay estimation (Andor, Frondel, & Vance, 2017; Champ, Moore, & Bishop, 2009; Mohammed, 2012; Ready, Champ, & Lawton, 2010). Namely, participants can revise their choices before continuing the survey:

*If we contact the restaurant now to place this order for you, will you be happy for us to proceed? [a) Yes, please place this order for me, b) No, I would like to change my choice]*

Then, we asked them if they would go to a restaurant offering similar food items. Answers are reported on a 5-Likert scale, ranging from “strongly agree” to “strongly disagree”. Revising one’s choice suggests low confidence in one’s preferences, increasing the risk of an intention-behaviour gap. Similarly, not wanting to go to a restaurant offering similar food items would suggest that participants would not make this choice in real life. Only 1.62% of the respondents have asked to revise their choices, and only 15.56% of them would either somewhat disagree or strongly disagree with going to a restaurant serving the same menus.

**Non-targeted pro-environmental behaviour:** At the end of the survey, we ask participants if they would like to donate an amount between £0 and £10 to a pro-environmental charity.<sup>14</sup> This task constitutes our non-targeted pro-environmental behaviour. We use pro-environmental donations as a proxy to measure respondents' willingness to engage in additional pro-environmental behaviours. Donations are consequential: we deduct them from the £20 voucher.

**Sample characteristics** We collected data from March 1<sup>st</sup> to April 24<sup>th</sup> of 2022. We pre-screened participants to select only native English speakers. We also excluded vegetarian and vegan participants. Attrition is low: only 4.1% of respondents were excluded for not finishing the survey. Table 22 in Appendix D shows descriptive statistics per treatment group. The median respondent is 35 years old, earns between £20,000 and £30,000 per year and has a Bachelor's degree. There is a good gender balance, with 49.9% of females, 49.2% of males and 0.9% of respondents considering themselves genderfluid or agender. Comparisons using the UK census data and the survey of personal income suggest that our sample is younger, slightly poorer and more educated than the UK population (see Figure 8 in Appendix D). Randomisation was successful. There are no significant differences across the four treatment groups regarding age, gender, income and education. About 98.28% of the *inference sample* has passed an attention check placed at the beginning of the survey,<sup>15</sup> and 99.75% of the screened sample has passed a focus check we placed after the pre-treatment questionnaire.<sup>16</sup> Furthermore, 81.69% of the

---

<sup>14</sup>Respondents are offered to give to the following charities: World Wide Fund (WWF), Friends of the Earth, Carbon Fund, Campaign against Climate Change, The Vegetarian Society, The Vegan Society, Extinction Rebellion, Woodland Trust. Alternatively, they can select "other" and write the name of their chosen charity.

<sup>15</sup>After consenting to participate to the survey experiment, respondents are screened based on whether they provide the correct answer to an attention check. Namely, they have to answer the following question on a 5-Likert scale, ranging from "not at all interested" to "extremely interested": "People are very busy these days, and many do not have time to follow what goes on in the government. We are testing whether people read questions. To show that you've read this much, answer both 'extremely interested' and 'very interested'."

<sup>16</sup>Participants have to answer the following question: "Most modern theories of decision making recognise that decisions do not take place in a vacuum. Individual preferences and knowledge, along with situational variables, can greatly impact the decision process. To demonstrate that you've read this much, just go ahead and select both red and green among the alternatives below. Based on the text you read above, what colour have you been asked to select?" They can select as many colours as they

participants passed a manipulation check between the food choice and the donation task.<sup>17</sup>

This suggests that respondents were attentive when taking the survey.

## IV.B Estimation Strategy

To estimate the side-effects of the policy on donations, we fit a linear model analogous to specifications (III.2) using ordinary least-square estimation (OLS):

$$Donation_i = \alpha^{SE} + \beta^{SE} Norm_i + \varepsilon_i^{SE} \quad (IV.1)$$

The outcome variable  $Donation_i$  is either a dummy equal to 1 when respondents choose to give and 0 otherwise, or a continuous variable ranging from 0 to 10 corresponding to the amount donated.  $Norm_i$  is the dummy capturing respondents' allocation to the social norm message.<sup>18</sup>

As we showed in Section II, the side-effects of policies can be decomposed into a behavioural spillover effect and a crowding in/out effect. A naive approach to disentangle these two effects consists of controlling for the targeted behaviour, here food choices, using a model analogous to specification (III.3), estimated with OLS:

$$Donation_i = \tilde{\alpha} + \tilde{\beta}^{BS} FoodChoice_i + \tilde{\beta}^C Norm_i + \tilde{\varepsilon}_i \quad (IV.2)$$

want from six colours. If they fail it, we show them the following message: *"The last question was here to check if you are being attentive. You did not answer it correctly. We are really interested in what you genuinely prefer. We kindly request you to read the questions more attentively."*

<sup>17</sup>This attention check was the following:

*Before being shown the restaurant menu, you were shown a message. What was the message about? [a) People changing diets to become climate-friendly, b) People changing their diets to lose weight, c) People changing their diets to respect animals' well-being, d) I was not shown any specific message, e) I do not remember any specific message displayed]*

<sup>18</sup>Therefore, estimate  $\hat{\beta}^{SE}$  corresponds to an intention-to-treat effect. In Appendix C, we assess the complier average causal effect by regressing  $Donation_i$  on a dummy equal to 1 when participants are shown the social norm message and correctly remember it in the manipulation check, and 0 otherwise. We instrument this dummy by respondents' random allocation to the social norm message.

The independent variable  $FoodChoice_i$  is either a dummy equal to 1 if respondents choose a vegetarian item, 0 otherwise, or a continuous variable capturing the carbon footprint of participants' food choices. As in Section III, coefficients  $\tilde{\beta}^{BS}$  and  $\tilde{\beta}^C$ , capture respectively the behavioural spillover effect and the crowding in/out effect of the social norm message. To tackle potential omitted variable biases, we instrument respondents' food choices by  $Menu_i$ , the dummy capturing allocation to the vegetarian chef's suggestion menu. We use a specification analogous to model (III.4),<sup>19</sup> estimated with two-stage least squares (2SLS):

$$\begin{aligned}
 1^{st} \text{ stage: } FoodChoice_i &= \alpha + \beta_1 Menu_i + \beta_2 Norm_i + \varepsilon_i \\
 2^{nd} \text{ stage: } Donation_i &= \alpha' + \beta^{BS} \widehat{FoodChoice}_i + \beta^C Norm_i + \varepsilon'_i
 \end{aligned}
 \tag{IV.3}$$

Using OLS, we estimate the policy's main effect on food choices by fitting the first stage of model (IV.3). We use probability linear models whenever the endogenous and outcome variables are binary. We relax the linearity assumption in robustness checks.<sup>20</sup> For all specifications, we also add lasso-selected controls to increase the precision of our estimates (see Appendix C, Belloni, Chernozhukov, and Hansen 2014). When presenting our results, we report standard p-values corrected for the false discovery rate (Benjamini & Hochberg, 1995), and p-values computed by re-randomising treatment allocation *à la* Young (2019). The latter approach is used as an extra robustness check to ensure leverage does not drive statistical significance. Finally, we also report adjusted confidence intervals for coefficient  $\beta^{BS}$  using Lee, McCrary, Moreira, and Porter (2022)'s procedure as confidence intervals derived from simple t-tests are likely to be biased with 2SLS estimation.

<sup>19</sup>In Appendix C, we test Assumption 1 by interacting food choices with respondents' exposure to the social norm nudge.

<sup>20</sup>We use probit models for specification (IV.1) and when checking for the robustness of the main effect of the social norm on food choices with the first stage of specification (IV.3). As an alternative to 2SLS estimation for specification (IV.3), we apply Rivers and Vuong (1988)'s two-step approach and a maximum likelihood estimation approach, as in Evans and Schwab (1995).

Table 2: ATE of the social norm message

Outcome	Chose vegetarian food (binary)	Food choice in kgCO2-eq
Specification	First stage	
Baseline	0.135*** (0.012)	23.400*** (0.871)
Social norm	0.067*** (0.016)	-2.751** (0.928)
	q<0.001	q=0.003
Plant-intensive chef's selection	0.115*** (0.016)	-7.875*** (0.928)
	q<0.001	q<0.001
Num.Obs.	2775	2775
R2	0.025	0.028

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: This table displays the effect of the social norm message and of being presented with a selection containing a majority of vegetarian items on respondents' decision to choose a vegetarian food item (first column) and on the carbon footprint of their choices (second column). Coefficients are estimated using OLS. Robust standard errors are displayed in parentheses. We apply [Benjamini and Hochberg \(1995\)](#) correction to conventional p-values (p). P-values of randomisation tests with 5,000 re-sampling are displayed last (q).

## IV.C Results

**Main effect:** Table 2 contains the effect of the social norm message on food choices.<sup>21</sup> The social norm nudge increases intention to choose vegetarian food by 6.7 percentage points and reduces the carbon footprint of food choices by 11.8%. Results are robust to non-linear probit specifications (see Table 8 in Appendix C).

**Side effects:** Table 3 displays the spillover effects of the social norm message on the binary decision to donate (Panel A) and the amount donated (Panel B). The first column contains the results of specification (IV.1) where donations are regressed on exposure to the social norm. This coefficient corresponds to the total side effects of the nudge. In both panels, these side effects are not significantly different from zero.

The second column displays the results obtained from fitting specification (IV.2). It corre-

<sup>21</sup>Analyses were conducted on R using the package *estimatr* ([Blair, Cooper, Coppock, Humphreys, & Sonnet, 2022](#)).



sponds to the naive approach for disentangling the crowding in/out from the behavioural spillover effects. The third column displays the results obtained from the two-stage least square regression (IV.3), where we instrument food choices by the default menu allocation. The effect of the social norm nudge on donations when controlling for food choices is not significantly different from zero, whether or not we instrument food choices. As such, we do not find evidence of crowding in/out effects.

However, we find suggestive evidence of a positive behavioural spillover effect. The correlation between food choices and donations is statistically significant (column two, specification (IV.2)). When instrumenting food choices, we find that individuals choosing a vegetarian dish are more likely to give by about 36 percentage points. There is no statistically significant effect on the amount donated after applying p-value correction. We do not observe a statistically significant difference between the instrumented and non-instrumented coefficients. The sign and magnitude of our estimates are robust when adding controls and when using non-linear specifications (see Tables 9 and 10 in Appendix C). P-values of randomisation tests indicate that outliers do not drive statistical significance.

**Strength of the IV:** Our instrumental variable should be relevant and exogenous. Results in Table 2 show a strong and highly significant effect of being allocated to a plant-intensive default menu on respondents' likelihood of choosing vegetarian food. The F statistic of the IV is 53.400 in the binary case and 71.998 when looking at the carbon footprint of food choices. This F-statistic is robust to adding controls (59.432 in the binary case, 64.140 with carbon footprint). This suggests that our instrument is strong (Bound, Jaeger, & Baker, 1995; Staiger & Stock, 1997). Recent empirical evidence from Gärtner (2018), Van Gestel, Adriaanse, and De Ridder (2020) and Ortmann, Ryvkin, Wilkening, and Zhang (2023) indicate default nudges affect people's decisions unconsciously, confirming priors in the literature (e.g., see Hansen and Jespersen 2013; Thaler and Sunstein 2009). Therefore, It seems unlikely that the instrumental variable affects donations other than through food choices.



Table 3: Total side effects, behavioural spillovers and crowding-in/out effects

<b>Panel A</b>					
Decision to donate (binary)					
Baseline	0.477*** (0.013)	0.443*** (0.014)	0.408*** (0.035)	0.525*** (0.015)	0.578*** (0.049)
Social norm	0.008 (0.019)	-0.004 (0.019)	-0.016 (0.022)	0.002 (0.019)	-0.006 (0.020)
Food choice	q=0.661	q=0.849	q=0.473	q=0.941	q=0.773
		0.178*** (0.022)	0.357* (0.166)	-0.002*** (0.000)	-0.005* (0.002)
			[0.004; 0.709] q<0.001		[-0.010; -0.002] q<0.001
R2	0.000	0.022		0.015	
<b>Panel B</b>					
Amount donated (in £)					
Baseline	3.309*** (0.108)	3.023*** (0.111)	2.870*** (0.272)	3.695*** (0.124)	3.956*** (0.389)
Social norm	-0.009 (0.151)	-0.109 (0.150)	-0.163 (0.175)	-0.063 (0.151)	-0.100 (0.161)
Food choice	q=0.952	q=0.473	q=0.338	q=0.680	q=0.528
		1.490*** (0.187)	2.286 (1.309)	-0.020*** (0.003)	-0.033 (0.019)
			[-0.495; 5.066] q=0.002		[-0.072; 0.006] q=0.002
R2	0.000	0.024		0.015	
Food choice		Binary	Binary	kgCO2-eq	kgCO2-eq
Specification	OLS	OLS	2SLS	OLS	2SLS
Num.Obs.	2775	2775	2775	2775	2775

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: This table displays the effect of food choices and the social norm message on the decision to donate (Panel A) and on the amount donated (Panel B). The first column shows the overall side effect of the social norm nudge on donations. The other columns show estimates of behavioural spillovers (effect of food choices on donations) and the crowding-in/out effect of the social norm message. The second and the fourth columns show estimates of the social norm nudge and food choices on donations. The third and the fifth columns display the same estimates where, this time, food choices are instrumented by default menu allocation. Robust standard errors are displayed in parentheses. We apply [Benjamini and Hochberg \(1995\)](#) correction to conventional p-values (p). The brackets display confidence intervals adjusted with [Lee et al. \(2022\)](#)'s procedure. P-values of randomisation tests with 5,000 re-sampling are displayed last (q).

**Discussion:** Our results indicate that eating vegetarian increases respondents’ willingness to do extra climate-friendly actions as proxied by our donation task. It is, however, important to note that we only estimate a local average treatment effect. When compliers’ profiles differ too much from the rest of the sample, this can affect the external validity of our results. Following [Marbach and Hangartner \(2020\)](#), we find that, compared to the average of the sample, compliers agree more with the idea that acting against climate change is a moral duty, order food online less frequently and agree less with the idea that British food should be meat-based (see [Figures 9](#) in [Appendix D](#)). Another caveat regards the hypothetical nature of food choices. An experimenter’s demand effect can inflate the average treatment effect of the social norm message, and the fact that the first behaviour is intentional could have induced participants willing to eat vegetarian to make a higher donation because they could not realise their intentions. Nevertheless, the fact that respondents who chose a vegetarian item report having exerted more effort for the environment in the post-treatment questionnaire seems to contradict this interpretation (see [Table 11](#) in [Appendix C](#)). Besides, our results align with evidence from field experiments showing positive behavioural spillover effects between pro-environmental actions ([Alacevich et al., 2021](#); [Comin & Rode, 2023](#)). Finally, our empirical strategy, as laid out in [Section III](#), partly relied on [Assumption 1](#). We fail to reject this assumption: the interaction between food choices and respondents’ exposure to the policy is not significantly different from zero (see [Table 12](#) in [Appendix C](#)).

As highlighted in [Section III](#), average treatment effects can hide heterogeneity. Respondents could perceive the social norm message differently depending on whether they are already conforming. We explore the heterogeneity of our causal effects in subsection [IV.D](#).

## **IV.D Heterogeneity Analysis**

How people perceive the social norm message might depend on how much they are willing to follow the norm. For instance, informing people of an upward trend in vegetarian diet adoption may trigger ditherers to change their behaviours, induce convinced meat-eaters to reaffirm

their preferences and be ignored by vegetarians having already exhausted all possibilities for improvement. In other words, the same social norm nudge will likely play on different cognitive processes for different profiles. We investigate this heterogeneity by classifying people into different profiles as part of an exploratory analysis.

**Training procedure:** In a separate survey, we showed 2,782 additional respondents the social norm message and then asked the following question:<sup>22</sup>

*Are you trying to change your diet to become more climate-friendly as well?*

- a) *No, I am not trying now, and I do not intend to try in future*
- b) *No, I am not trying now, but I might consider changing my diet to be more-climate-friendly in future*
- c) *Yes, I am trying to change my diet now to become more climate-friendly*
- d) *Yes, I have already changed my diet to be more climate-friendly*

We assume that asking this question right after respondents see the social norm message reveals their willingness to follow the norm. This question allows us to identify four types: the *transitioned type*, who is already conforming with the social norm; the *trying type*, who is willing to conform; the *hesitant type*, who considers doing so in the future, and the *unwilling type* who does not want to conform. We use respondents' answers to train a gradient tree boosting machine learning classifier (GBM) to predict the types of respondents in the *inference sample*.<sup>23</sup>

As with Random Forest, GBM fits multiple decision trees. Here, the difference is that each

---

<sup>22</sup>This question is part of another treatment arm designed for another research project testing if inducing people to think about their choices increases the effectiveness of social norm nudges. See [Banerjee and Picard \(2023\)](#) for more details.

<sup>23</sup>Despite having excluded vegan and vegetarian participants, 12,6% of respondents chose the last answer. We see three explanations for this apparent contradiction. First, the screening was based on social demographic information gathered by Prolific, our data provider. As such, people may have changed their diets between when they answered the Prolific questionnaire and when they took our survey. Second, answers can also capture intentions rather than behaviours. Third, the phrasing of this answer could have been perceived as vague enough to allow non-vegetarian participants to select it without contradicting their actual behaviour.

additional decision tree is fitted on the errors made by the previous one (Friedman, 2001). We explain the algorithm in detail in Appendix B. To test the robustness of our predictions, we train five additional classification algorithms: random forest, a multinomial regression model, an ordered logit model, linear discriminant analysis, and quadratic discriminant analysis.<sup>24</sup> We estimate the average performance of GBM using nested 10×10 folds cross-validation. Overall, GBM performs twice better than chance. Readers can refer to Appendix B for an extended discussion on the estimation of performances and an analysis of the predictive power of each predictor. The four classes predicted by GBM seem very similar to their counterparts in the training set (see density plots 10, 11 and 12 in Appendix D).

**Profile of predicted types:** Table 24 in Appendix C displays the extent to which each type differs from the average for each covariate. Respondents predicted to be *unwilling* to change their diet to follow the norm agree less with the idea that acting against climate change is a moral duty and agree more with the idea that climate change is exaggerated compared to the average. They also know less about the environmental impact of food. *Unwilling* respondents are older, less educated, less likely to live in London and more likely to be male and conservative than the average. Respondents in this group tend to agree more with the idea that typical British food should be meat-based and report a stronger preference for meat-based food. They are also less likely to follow a specific diet.

Respondents who are predicted to be *hesitant* about reducing meat to follow the norm are slightly more likely to live in an area with a higher unemployment rate and a lower number of students. Their area of residence is also less likely to be rural than the average. These respondents agree less with the idea that acting against climate change is a moral duty, know less about the environmental impact of food and are less confident in their knowledge of the environmental impact of food. They are younger, less educated, slightly more likely to be female, poorer, and more likely to live in the same area than their area of birth. They also agree

---

<sup>24</sup>For more information on how these algorithms work, the reader can refer to Gareth, Daniela, Trevor, and Robert (2013).

more with the idea that British food should be meat-based, report a stronger preference for meat-based food and order food online more frequently than the average.

Respondents predicted to be *trying* to follow the norm are slightly more likely to live in an area where the unemployment rate is lower and the number of students is higher than the average. They agree more with the idea that acting against climate change is a moral duty, agree less that climate change is exaggerated, know more about the environmental impact of food and are more confident in their knowledge. Respondents in this group are older, more educated, more likely to have moved out of their area of birth, more likely to live in London, richer and less conservative than the average. They also report a lower preference for meat-based food, are less likely to follow a specific diet and order food online less frequently than the average.

Finally, respondents predicted to have *transitioned* to vegetarian diets are slightly more likely to live in a rural area with a lower share of unemployment. They agree more with the idea that acting against climate change is a moral duty, agree less with the idea that climate change is exaggerated, have a better knowledge of the environmental impact of food and are more confident in their knowledge. Respondents in this group are more educated, more likely to be female, to have moved out of their area of birth, and less conservative than the average. They agree less that British food should be meat-based, report a lower preference for meat-based food, and are likelier to follow a specific diet. They also order food online less frequently than the average of the sample.

In what follows, we estimate the main effect of the social norm message and its crowding-out/in effect for each predicted profile.

**Identification strategy:** When looking at the average treatment effect of the social norm nudge for each predicted type, we use the *unwilling* type as our reference group and fit the

following nested probability linear model:<sup>25</sup>

$$FoodChoice_i = \alpha + \sum_{k \in \Omega_-} \mathbf{1}_k \delta_k + \sum_{k \in \Omega} \mathbf{1}_k \beta_k Norm_i + u_i \quad (IV.4)$$

$$\Omega_- = \{\text{hesitant, trying, transitioned}\}$$

$$\Omega = \{\text{unwilling, hesitant, trying, transitioned}\}$$

And:

$$\mathbf{1}_k = \begin{cases} 1, & \text{if individual } i \text{ type } k \\ 0, & \text{otherwise} \end{cases}$$

Coefficient  $\beta_k$  is the average effect of the social norm nudge conditional on being predicted to be of type  $k$ . When looking at the total spillover effect of the social norm message conditional on respondents' predicted types, we fit the following model:

$$Donation_i = \alpha + \sum_{k \in \Omega_-} \mathbf{1}_k \delta_k + \sum_{k \in \Omega} \mathbf{1}_k \beta_k Norm_i + u_i \quad (IV.5)$$

Here again,  $\beta_k$  is the average side-effect of the social norm nudge conditional on being predicted to be of type  $k$ . To investigate heterogeneity in the crowding-in/out effect, we fit the following model:

$$Donation_i = \alpha + \sum_{k \in \Omega_-} \mathbf{1}_k \delta_k + \sum_{k \in \Omega} \mathbf{1}_k \beta_k Norm_i + \beta_2 \widehat{FoodChoice}_i + \varepsilon_i \quad (IV.6)$$

Here, the coefficient  $\beta_k$  is the crowding in/out effect of the social norm message conditional of being predicted to be of type  $k$ .  $\widehat{FoodChoice}_i$  captures respondents' food choices instrumented by the default menu allocation. In testing the robustness of our results, we also fit these models

---

<sup>25</sup>Such a specification is equivalent to fitting four separate linear models for each predicted profile.

with the predictions yielded by five other classification algorithms and the predictions yielded by over-sampling the *unwilling* and *transitioned* categories that contain fewer observations. We also re-estimate our GBM model by including previously excluded predictors found to contain too many missing values: income and political beliefs. We compute the p-values from re-randomisation tests to ensure outliers do not drive statistical significance (Young, 2019).

**Results:** Table 4 displays the results obtained by fitting equation (IV.4). Being predicted to be aligned with the social norm message is positively associated with a higher likelihood of choosing a vegetarian dish and negatively associated with the carbon footprint of food choices. After adjusting p-values, the social norm message is only effective on the predicted *trying* type at the extensive margin (+10.5 percentage points to choose a vegetarian meal) and in reducing the emissions of the predicted *hesitant* type (−18 pp of GHG emissions). These results are robust when looking at those obtained using the predictions of the other algorithms (see Tables 13 and 14 in Appendix C). Coefficients are of the same sign across all the algorithms and globally of the same order of magnitude. P-values of re-randomisation tests confirm that leverage does not drive statistical significance.

Table 5 shows the results of regression (IV.5) in the first two columns and regression (IV.6) in the last four columns. Although not significant after p-value correction, the social norm message yields negative total side effects on the amount given for the predicted *trying type*. We see that the social norm message drives this effect when controlling whether participants chose a vegetarian dish instrumented by default menu allocation. The social norm message crowds out the amount they donate by about £0.829 and the likelihood of donating by 9.3 percentage points. This crowding out effect is globally robust when looking at the predictions generated by the other algorithms (see Tables 17, 18, 19 and 20 in Appendix C). Again, p-values of re-randomisation tests confirm that leverage is not driving statistical significance. We also observe suggestive evidence of a crowding-in effect among the predicted *unwilling*. However, this effect is not significant after correcting for multiple hypothesis testing.

Table 4: Main effect of the social norm message conditional on respondents' types

Specification	Nested OLS model	
	Chose vegetarian food	Food choice in kgCO2-eq
Unwilling (baseline)	0.091*** (0.023)	24.412*** (2.122)
Hesitant	0.071*** (0.026)	-2.960 (2.341)
Trying	0.159*** (0.031)	-9.205*** (2.401)
Transitioned	0.337*** (0.064)	-13.954*** (3.086)
Social norm × Unwilling	-0.025 (0.030) q=0.407	3.076 (3.138) q=0.326
Social norm × Hesitant	0.039 (0.020)	-3.851** (1.328)
Social norm × Trying	0.108*** (0.033) q=0.051	-2.261 (1.577) q=0.006
Social norm × Transitioned	0.148 (0.077) q=0.001	-1.426 (2.899) q=0.152
	q=0.056	q=0.621
Num.Obs.	2730	2730
R2	0.068	0.032

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: This table displays the main effect of the social norm message on the likelihood of choosing vegetarian food (first column) and on the carbon footprint of food choices (second column) conditional on respondents' predicted types. For instance, coefficients labelled "Social norm × Trying" correspond to the average effect of the social norm on the predicted *trying* (the difference between control units and treatment units in this subsample). Coefficients labelled *Trying* correspond to the difference between the control units in the *trying* sample with the control units in the *unwilling* sample, the baseline. Robust standard errors are displayed in parentheses. We apply [Benjamini and Hochberg \(1995\)](#) correction to conventional p-values (p). P-values of randomisation tests with 5,000 re-sampling are displayed last (q).



Table 5: Crowding-in/out effects conditional on predicted types.

Specification	Nested OLS model		Nested 2SLS model			
	Amount (in £)	Decision (binary)	Amount (in £)	Decision (binary)	Amount (in £)	Decision (binary)
Food choice			Chose vegetarian food		Food choice in kgCO2-eq	
Unwilling (baseline)	1.348*** (0.234)	0.207*** (0.032)	1.168*** (0.259)	0.180*** (0.035)	2.060*** (0.509)	0.314*** (0.065)
Food choice			1.963 (1.225)	0.295* (0.154)	-0.029 (0.018)	-0.004* (0.002)
			q=0.003	q=0.001	q=0.002	q=0.001
Hesitant	1.552*** (0.273)	0.233*** (0.037)	1.413*** (0.285)	0.212*** (0.039)	1.466*** (0.282)	0.220*** (0.039)
Trying	3.289*** (0.316)	0.419*** (0.040)	2.977*** (0.371)	0.372*** (0.047)	3.021*** (0.361)	0.379*** (0.046)
Transitioned	3.538*** (0.556)	0.436*** (0.066)	2.876*** (0.723)	0.336*** (0.087)	3.131*** (0.617)	0.374*** (0.073)
Social norm × Unwilling	0.699 (0.377)	0.073 (0.049)	0.748 (0.380)	0.080 (0.049)	0.789 (0.390)	0.086 (0.051)
	q=0.064	q=0.134	q=0.049	q=0.107	q=0.040	q=0.085
Social norm × Hesitant	0.135 (0.197)	0.027 (0.026)	0.059 (0.204)	0.016 (0.027)	0.023 (0.212)	0.010 (0.027)
	q=0.487	q=0.299	q=0.778	q=0.544	q=0.915	q=0.706
Social norm × Trying	-0.664* (0.304)	-0.070 (0.036)	-0.877** (0.331)	-0.102** (0.039)	-0.730* (0.307)	-0.080* (0.036)
	q=0.032	q=0.052	q=0.008	q=0.011	q=0.020	q=0.028
Social norm × Transitioned	-0.366 (0.652)	0.001 (0.075)	-0.658 (0.691)	-0.042 (0.081)	-0.408 (0.644)	-0.005 (0.073)
	q=0.554	q=0.992	q=0.336	q=0.595	q=0.520	q=0.946
Num.Obs.	2730	2730	2730	2730	2730	2730
R2	0.051	0.051				

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: This table displays the effect of the social norm message on the decision to donate (columns 2, 4 and 6) and on the amount donated (columns 1, 3, and 5) conditional on the predicted type. The first two column shows the overall side effect of the social norm nudge on donations. The crowding-in/out effect of the social norm message is then estimated in the other columns, controlling for food choices instrumented by default menu allocation. For instance, coefficients labelled "Social norm × Trying" correspond to the effect of the social norm on the predicted *trying*. Coefficients labelled *Trying* correspond to the difference between the control units in the *trying* sample with the control units in the *unwilling* sample, the baseline. Robust standard errors are displayed in parentheses. We apply [Benjamini and Hochberg \(1995\)](#) correction to conventional p-values (p). P-values of randomisation tests with 5,000 re-sampling are displayed last (q).

**Discussion:** Our results indicate that the social norm message is effective on those who consider making changes (the *hesitant* type) and those making efforts (the *trying* type). For the *hesitant* type, the message induces a decrease in the carbon footprint of food choices but no increase in the uptake of vegetarian options after correcting p-values. This apparent paradox might be caused by the predicted *hesitants* switching from choosing carbon-intensive meat to less intensive meat options, resulting in decreased carbon emission but no statistically significant increase in vegetarian food uptake. Conversely, the message induces an increase in the uptake of vegetarian food but no significant decrease in carbon emissions among the predicted *trying*. This might be because participants classed as *trying* switch from less intensive meat options to vegetarian options, implying no statistically significant decrease in carbon emissions. Furthermore, we cannot reject the hypothesis that the message does not affect respondents predicted to be *unwilling*. Although the absence of evidence is not evidence of the absence, this null result supports a common assumption in the literature that nudges are ineffective for those unwilling to change (Thaler & Sunstein, 2009). Similarly, the message does not significantly alter the choices of the predicted *transitioned* type after correcting p-values. The *transitioned* respondents have the highest share of controlled units choosing vegetarian food. As such, it can be that *transitioned* respondents have already exhausted all possibilities of improvement.

We find clear and robust evidence that the social norm message crowds out donations of the predicted *trying* type. Our model in Section II suggests the *trying* respondents may have treated the social norm message as an extrinsic pressure to choose vegetarian food. Following Deci and Ryan (2010), this would have induced them to slacken once the extrinsic pressure vanishes. Truelove, Carrico, Weber, Raimi, and Vandenberg (2014)'s theoretical framework provides a similar interpretation. For the authors, policies can induce people to act to repair a morally threatened identity, inducing moral licensing once the identity is repaired. Interestingly, the social norm message does not produce a similar crowding-out effect for the predicted *hesitants*. This can be because respondents classed as *trying* are more aware of the environmen-

tal impact of diets, making them more prone to guilt when exposed to our message. We also find suggestive evidence of a crowding-in effect of the social norm message on the predicted *unwilling*. Although not statistically significant, this would explain why the crowding-in/out effect is close to zero when aggregating the four categories. The fact that the predicted *unwilling* did not alter their food choices but chose to contribute more through donations suggests that this subsample may have engaged in moral cleansing (Sachdeva, Iliev, & Medin, 2009). Moral cleansing describes pro-social acts undertaken to repair deprecated moral self-worth. However, this interpretation should be considered with caution, given the fragility of this result.

## V Conclusion

In Section II, we model the side effects of behavioural policies as the sum of two effects. The first effect, commonly referred to as a *behavioural spillover* (Thøgersen, 1999), emerges when a policy successfully shifts a targeted decision, subsequently encouraging or discouraging further pro-environmental decisions. We label the second a *crowding-in/out* effect. It captures the policy’s impact on our propensity to engage further. Its sign, therefore, depends on the nature of the policy used. We use an instrumental variable to dissociate the behavioural spillovers from the crowding-in/out effects of the social norm message. Furthermore, we explore heterogeneity in the effects of the social norm message by identifying profiles expected to respond differently to the message using an additional treatment arm. We then use this extra survey data to classify respondents in the initial sample into different types using machine learning.

Our results are consistent with other studies that use an instrumental variable approach to estimate behavioural spillovers between pro-environmental behaviours. Comin and Rode (2023) find evidence that installing solar panels increase support for pro-environmental policies. Similarly, Alacevich et al. (2021) find that sorting waste results in an overall decrease in the amount of waste generated. Our paper finds that intentions to choose vegetarian options

foster pro-environmental donations. As such, on top of yielding large reductions in greenhouse gas emissions (Green et al., 2015; Riahi et al., 2022), cutting on meat seems to increase people's willingness to do more.

In this regard, using social norm messaging to promote vegetarianism is an effective strategy for people that we predict as trying to change their diets to follow the norm and those hesitating about doing so. However, we only observe a significant reduction in the carbon footprint of food choices for the predicted *hesitants*. Besides, the message crowds out the predicted *tryings'* donations. This crowding-out effect outweighs the positive behavioural spillover effect. We do not observe a similar crowding-out effect on the respondents predicted to be hesitant. This suggests that policymakers seeking to use social norm messages to reduce the environmental impact of food choices should target this population segment.

When it comes to increasing the uptake of vegetarian choices, our experimental findings indicate no "free lunch": when the social norm message appears effective in fostering vegetarian food choices, it comes at the cost of crowding out further engagement. This calls for more empirical evidence on whether other types of policies yield similar results. We have laid out a theoretical and empirical framework to assess these questions, dissociating the side effects triggered by policy interventions and exploring heterogeneity in their effects. We hope this paper will provide the methodological foundations for further research.

### **Informed consent and ethics approval**

All participants participated in the study with their informed consent. This study aligned with the London School of Economics and Political Science research ethics guidelines. The study was approved vide reference 38224.

### **Funding Statement**

This study was funded by the Royal Geographic Society (RGS-IBG) Frederick Soddy Postgraduate Award vide reference FSPA 05.21.

## Competing interests

The authors have no competing interests to declare.

## Code and Data Availability

The code and data for the analysis are available upon request.

## References

- Abrahamse, W., & Steg, L. (2013). Social influence approaches to encourage resource conservation: A meta-analysis. *Global environmental change, 23*(6), 1773–1785.
- Akerlof, G. A., & Kranton, R. E. (2000). Economics and identity. *The quarterly journal of economics, 115*(3), 715–753.
- Alacevich, C., Bonev, P., & Söderberg, M. (2021). Pro-environmental interventions and behavioral spillovers: Evidence from organic waste sorting in sweden. *Journal of Environmental Economics and Management, 108*, 102470.
- Allcott, H. (2011). Social norms and energy conservation. *Journal of public Economics, 95*(9-10), 1082–1095.
- Alt, M., Bruns, H., & DellaValle, N. (2023). The more the better?-synergies of prosocial interventions and effects on behavioral spillovers. *Synergies of prosocial interventions and effects on behavioral spillovers (June 27, 2023)*.
- Alt, M., & Gallier, C. (2022). Incentives and intertemporal behavioral spillovers: A two-period experiment on charitable giving. *Journal of Economic Behavior & Organization, 200*, 959–972.
- Andersson, M., & von Borgstede, C. (2010). Differentiation of determinants of low-cost and high-cost recycling. *Journal of Environmental Psychology, 30*(4), 402–408.
- Andor, M. A., & Fels, K. M. (2018). Behavioral economics and energy conservation—a systematic review of non-price interventions and their causal effects. *Ecological economics, 148*, 178–210.
- Andor, M. A., Frondel, M., & Vance, C. (2017). Mitigating hypothetical bias: Evidence on the effects of correctives from a large field study. *Environmental and Resource Economics, 68*(3), 777–796.
- Athey, S., & Imbens, G. W. (2015). Machine learning methods for estimating heterogeneous causal effects. *stat, 1050*(5), 1–26.
- Baca-Motes, K., Brown, A., Gneezy, A., Keenan, E. A., & Nelson, L. D. (2013). Commitment and behavior change: Evidence from the field. *Journal of Consumer Research, 39*(5), 1070–1084.
- Banerjee, S., Galizzi, M. M., John, P., & Mourato, S. (2022). What works best in promoting climate citizenship? a randomised, systematic evaluation of nudge, think,

- boost and nudge+. *A Randomised, Systematic Evaluation of Nudge, Think, Boost and Nudge+* (April 26, 2022). *Research Square pre-prints*.
- Banerjee, S., & Picard, J. (2023). Thinking through norms can make them more effective. experimental evidence on reflective climate policies in the uk. *Journal of Behavioral and Experimental Economics*, 102024.
- Belloni, A., Chernozhukov, V., & Hansen, C. (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2), 608–650.
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1), 289–300.
- Bicchieri, C. (2016). *Norms in the wild: How to diagnose, measure, and change social norms*. Oxford University Press.
- Blair, G., Cooper, J., Coppock, A., Humphreys, M., & Sonnet, L. (2022). *estimatr: Fast estimators for design-based inference* [Computer software manual]. (<https://declaredesign.org/r/estimatr/>, <https://github.com/DeclareDesign/estimatr>)
- Blondin, S., Attwood, S., Vennard, D., & Mayneris, V. (2022). Environmental messages promote plant-based food choices: An online restaurant menu study. *World Resources Institute*.
- Bonev, P. (2023). Behavioral spillovers.
- Bound, J., Jaeger, D. A., & Baker, R. M. (1995). Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *Journal of the American statistical association*, 90(430), 443–450.
- Bratt, C. (1999). The impact of norms and assumed consequences on recycling behavior. *Environment and behavior*, 31(5), 630–656.
- Bryan, C. J., Tipton, E., & Yeager, D. S. (2021). Behavioural science is unlikely to change the world without a heterogeneity revolution. *Nature human behaviour*, 5(8), 980–989.
- Carrico, A. R., Raimi, K. T., Truelove, H. B., & Eby, B. (2018). Putting your money where your mouth is: an experimental test of pro-environmental spillover from reducing meat consumption to monetary donations. *Environment and Behavior*, 50(7), 723–748.
- Carrico, A. R., & Riemer, M. (2011). Motivating energy conservation in the workplace: An evaluation of the use of group-level feedback and peer education. *Journal of environmental psychology*, 31(1), 1–13.
- Champ, P. A., Moore, R., & Bishop, R. C. (2009). A comparison of approaches to mitigate hypothetical bias. *Agricultural and Resource Economics Review*, 38(2), 166–180.



- Cialdini, R. B., & Jacobson, R. P. (2021). Influences of social norms on climate change-related behaviors. *Current Opinion in Behavioral Sciences*, 42, 1–8.
- Clot, S., Della Giusta, M., & Jewell, S. (2022). Once good, always good? testing nudge’s spillovers on pro environmental behavior. *Environment and Behavior*, 54(3), 655–669.
- Comin, D. A., & Rode, J. (2023). *Do green users become green voters?* (Tech. Rep.). National Bureau of Economic Research.
- Costa, D. L., & Kahn, M. E. (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.
- Deci, E. L., & Ryan, R. M. (2010). Intrinsic motivation. *The corsini encyclopedia of psychology*, 1–2.
- Dolan, P., & Galizzi, M. M. (2014). Because i’m worth it: a lab-field experiment on the spillover effects of incentives in health.
- Evans, W. N., & Schwab, R. M. (1995). Finishing high school and starting college: Do catholic schools make a difference? *The Quarterly Journal of Economics*, 110(4), 941–974.
- Farrow, K., Grolleau, G., & Ibanez, L. (2017). Social norms and pro-environmental behavior: A review of the evidence. *Ecological Economics*, 140, 1–13.
- Ferraro, P. J., Miranda, J. J., & Price, M. K. (2011). The persistence of treatment effects with norm-based policy instruments: evidence from a randomized environmental policy experiment. *American Economic Review*, 101(3), 318–22.
- Festinger, L. (1962). Cognitive dissonance. *Scientific American*, 207(4), 93–106.
- Fornara, F., Carrus, G., Passafaro, P., & Bonnes, M. (2011). Distinguishing the sources of normative influence on proenvironmental behaviors: The role of local norms in household waste recycling. *Group Processes & Intergroup Relations*, 14(5), 623–635.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189–1232.
- Gareth, J., Daniela, W., Trevor, H., & Robert, T. (2013). *An introduction to statistical learning: with applications in r*. Springer.
- Gärtner, M. (2018). The prosociality of intuitive decisions depends on the status quo. *Journal of Behavioral and Experimental Economics*, 74, 127–138.
- Geiger, S. J., Brick, C., Nalborczyk, L., Bosshard, A., & Jostmann, N. B. (2021). More green than gray? toward a sustainable overview of environmental spillover effects: A bayesian meta-analysis. *Journal of Environmental Psychology*, 78, 101694.
- Gneezy, A., Imas, A., Brown, A., Nelson, L. D., & Norton, M. I. (2012). Paying to be nice: Consistency and costly prosocial behavior. *Management Science*, 58(1), 179–187.
- Goldstein, N. J., Cialdini, R. B., & Griskevicius, V. (2008). A room with a viewpoint: Using social norms to motivate environmental conservation in hotels. *Journal of*

- consumer Research*, 35(3), 472–482.
- Green, R., Milner, J., Dangour, A. D., Haines, A., Chalabi, Z., Markandya, A., ... Wilkinson, P. (2015). The potential to reduce greenhouse gas emissions in the uk through healthy and realistic dietary change. *Climatic Change*, 129(1), 253–265.
- Handgraaf, M. J., De Jeude, M. A. V. L., & Appelt, K. C. (2013). Public praise vs. private pay: Effects of rewards on energy conservation in the workplace. *Ecological Economics*, 86, 86–92.
- Hansen, P. G., & Jespersen, A. M. (2013). Nudge and the manipulation of choice: A framework for the responsible use of the nudge approach to behaviour change in public policy. *European Journal of Risk Regulation*, 4(1), 3–28.
- Hartmann, P., Marcos, A., & Barrutia, J. M. (2023). Carbon tax salience counteracts price effects through moral licensing. *Global Environmental Change*, 78, 102635.
- Kristofferson, K., White, K., & Peloza, J. (2014). The nature of slacktivism: How the social observability of an initial act of token support affects subsequent prosocial action. *Journal of Consumer Research*, 40(6), 1149–1166.
- Künzel, S. R., Sekhon, J. S., Bickel, P. J., & Yu, B. (2019). Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the national academy of sciences*, 116(10), 4156–4165.
- Lacasse, K. (2016). Don't be satisfied, identify! strengthening positive spillover by connecting pro-environmental behaviors to an "environmentalist" label. *Journal of Environmental Psychology*, 48, 149–158.
- Lapinski, M. K., Rimal, R. N., DeVries, R., & Lee, E. L. (2007). The role of group orientation and descriptive norms on water conservation attitudes and behaviors. *Health communication*, 22(2), 133–142.
- Lee, D. S., McCrary, J., Moreira, M. J., & Porter, J. (2022). Valid t-ratio inference for iv. *American Economic Review*, 112(10), 3260–3290.
- Liu, Y., Kua, H., & Lu, Y. (2021). Spillover effects from energy conservation goal-setting: A field intervention study. *Resources, Conservation and Recycling*, 170, 105570.
- Löfgren, Å., & Nordblom, K. (2020). A theoretical framework of decision making explaining the mechanisms of nudging. *Journal of Economic Behavior & Organization*, 174, 1–12.
- Lokhorst, A. M., Werner, C., Staats, H., van Dijk, E., & Gale, J. L. (2013). Commitment and behavior change: A meta-analysis and critical review of commitment-making strategies in environmental research. *Environment and behavior*, 45(1), 3–34.
- Maki, A., Carrico, A. R., Raimi, K. T., Truelove, H. B., Araujo, B., & Yeung, K. L. (2019). Meta-analysis of pro-environmental behaviour spillover. *Nature Sustainability*, 2(4), 307–315.
- Marbach, M., & Hangartner, D. (2020). Profiling compliers and noncompliers for instrumental-variable analysis. *Political Analysis*, 28(3), 435–444.



- Margetts, E. A., & Kashima, Y. (2017). Spillover between pro-environmental behaviours: The role of resources and perceived similarity. *Journal of Environmental Psychology, 49*, 30–42.
- Mazar, N., & Zhong, C.-B. (2010). Do green products make us better people? *Psychological science, 21*(4), 494–498.
- Melnyk, V., van Herpen, E., Trijp, H., et al. (2010). The influence of social norms in consumer decision making: A meta-analysis. *ACR North American Advances*.
- Merritt, A. C., Effron, D. A., & Monin, B. (2010). Moral self-licensing: When being good frees us to be bad. *Social and personality psychology compass, 4*(5), 344–357.
- Mohammed, E. Y. (2012). Contingent valuation responses and hypothetical bias: mitigation effects of certainty question, cheap talk, and pledging. *Environmental Economics*(3, Iss. 3), 62–71.
- Nigbur, D., Lyons, E., & Uzzell, D. (2010). Attitudes, norms, identity and environmental behaviour: Using an expanded theory of planned behaviour to predict participation in a kerbside recycling programme. *British journal of social psychology, 49*(2), 259–284.
- Nolan, J. M., Schultz, P. W., Cialdini, R. B., Goldstein, N. J., & Griskevicius, V. (2008). Normative social influence is underdetected. *Personality and social psychology bulletin, 34*(7), 913–923.
- Ortmann, A., Ryvkin, D., Wilkening, T., & Zhang, J. (2023). Defaults and cognitive effort. *Journal of Economic Behavior & Organization, 212*, 1–19.
- Picard, J. (2021). Nudging virtuous behaviours without crowding-out other ones: micro-foundations to behavioural interventions. Available at SSRN 3784009.
- Ready, R. C., Champ, P. A., & Lawton, J. L. (2010). Using respondent uncertainty to mitigate hypothetical bias in a stated choice experiment. *Land Economics, 86*(2), 363–381.
- Reese, G., Loew, K., & Steffgen, G. (2014). A towel less: Social norms enhance pro-environmental behavior in hotels. *The Journal of Social Psychology, 154*(2), 97–100.
- Rhodes, N., Shulman, H. C., & McClaran, N. (2020). Changing norms: A meta-analytic integration of research on social norms appeals. *Human Communication Research, 46*(2-3), 161–191.
- Riahi, K., Schaeffer, R., Arango, J., Calvin, K., Guivarch, C., Hasegawa, T., ... others (2022). Mitigation pathways compatible with long-term goals. *IPCC, 2022: Climate Change 2022: Mitigation of Climate Change. Contribution of Working Group III to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change*.
- Richter, I., Thøgersen, J., & Klöckner, C. A. (2018). A social norms intervention going wrong: Boomerang effects from descriptive norms information. *Sustainability, 10*(8), 2848.
- Rivers, D., & Vuong, Q. H. (1988). Limited information estimators and exogeneity tests for simultaneous probit models. *Journal of econometrics, 39*(3), 347–366.

- Sachdeva, S., Iliev, R., & Medin, D. L. (2009). Sinning saints and saintly sinners: The paradox of moral self-regulation. *Psychological science*, 20(4), 523–528.
- Salmivaara, L., & Lankoski, L. (2019). Promoting sustainable consumer behaviour through the activation of injunctive social norms: A field experiment in 19 workplace restaurants. *Organization & Environment*, 1086026619831651.
- Scarborough, P., Appleby, P. N., Mizdrak, A., Briggs, A. D., Travis, R. C., Bradbury, K. E., & Key, T. J. (2014). Dietary greenhouse gas emissions of meat-eaters, fish-eaters, vegetarians and vegans in the uk. *Climatic change*, 125(2), 179–192.
- Schultz, W. P., Khazian, A. M., & Zaleski, A. C. (2008). Using normative social influence to promote conservation among hotel guests. *Social influence*, 3(1), 4–23.
- Shukla, Skea, Slade, Khouradajie, A., van Diemen, McCollum, . . . Malley (2022). Mitigation of climate change. contribution of working group iii to the sixth assessment report of the intergovernmental panel on climate change. *Cambridge University Press*.
- Smith, A. (1853). *The theory of moral sentiments*. HG Bohn.
- Sparkman, G., & Walton, G. M. (2017). Dynamic norms promote sustainable behavior, even if it is counternormative. *Psychological science*, 28(11), 1663–1674.
- Sparkman, G., Weitz, E., Robinson, T. N., Malhotra, N., & Walton, G. M. (2020). Developing a scalable dynamic norm menu-based intervention to reduce meat consumption. *Sustainability*, 12(6), 2453.
- Staiger, D., & Stock, J. H. (1997). Instrumental variables regression with weak instruments. *Econometrica*, 65(3), 557–586.
- Stea, S., & Pickering, G. J. (2019). Optimizing messaging to reduce red meat consumption. *Environmental Communication*, 13(5), 633–648.
- Steinhorst, J., & Matthies, E. (2016). Monetary or environmental appeals for saving electricity?—potentials for spillover on low carbon policy acceptability. *Energy Policy*, 93, 335–344.
- Stewart, C., Piernas, C., Cook, B., & Jebb, S. A. (2021). Trends in uk meat consumption: analysis of data from years 1–11 (2008–09 to 2018–19) of the national diet and nutrition survey rolling programme. *The Lancet Planetary Health*, 5(10), e699–e708.
- Testa, F., Russo, M. V., Cornwell, T. B., McDonald, A., & Reich, B. (2018). Social sustainability as buying local: effects of soft policy, meso-level actors, and social influences on purchase intentions. *Journal of Public Policy & Marketing*, 37(1), 152–166.
- Thaler, R. H., & Sunstein, C. R. (2009). *Nudge: Improving decisions about health, wealth, and happiness*. Penguin.
- Thøgersen, J. (1999). Spillover processes in the development of a sustainable consumption pattern. *Journal of economic psychology*, 20(1), 53–81.
- Tiefenbeck, V., Staake, T., Roth, K., & Sachs, O. (2013). For better or for worse? empirical evidence of moral licensing in a behavioral energy conservation campaign.

- Energy Policy*, 57, 160–171.
- Truelove, H. B., Carrico, A. R., Weber, E. U., Raimi, K. T., & Vandenberg, M. P. (2014). Positive and negative spillover of pro-environmental behavior: An integrative review and theoretical framework. *Global Environmental Change*, 29, 127–138.
- Van der Werff, E., Steg, L., & Keizer, K. (2013). It is a moral issue: The relationship between environmental self-identity, obligation-based intrinsic motivation and pro-environmental behaviour. *Global environmental change*, 23(5), 1258–1265.
- Van der Werff, E., Steg, L., & Keizer, K. (2014). I am what i am, by looking past the present: the influence of biospheric values and past behavior on environmental self-identity. *Environment and Behavior*, 46(5), 626–657.
- Van Gestel, L., Adriaanse, M., & De Ridder, D. (2020). Do nudges make use of automatic processing? unraveling the effects of a default nudge under type 1 and type 2 processing. *Comprehensive Results in Social Psychology*, 1–21.
- Van Rookhuijzen, M., De Vet, E., & Adriaanse, M. A. (2021). The effects of nudges: One-shot only? exploring the temporal spillover effects of a default nudge. *Frontiers in Psychology*, 12.
- Wenzig, J., & Gruchmann, T. (2018). Consumer preferences for local food: Testing an extended norm taxonomy. *Sustainability*, 10(5), 1313.
- Wolstenholme, E., Poortinga, W., & Whitmarsh, L. (2020). Two birds, one stone: The effectiveness of health and environmental messages to reduce meat consumption and encourage pro-environmental behavioral spillover. *Frontiers in psychology*, 11, 577111.
- Xu, L., Zhang, X., & Ling, M. (2018). Spillover effects of household waste separation policy on electricity consumption: evidence from hangzhou, china. *Resources, Conservation and Recycling*, 129, 219–231.
- Young, A. (2019). Channeling fisher: Randomization tests and the statistical insignificance of seemingly significant experimental results. *The Quarterly Journal of Economics*, 134(2), 557–598.

## A Appendix: Proofs

*Proof. (Main effect of policies)* The effect of a policy altering  $\theta_1 \in \{\eta_1, \epsilon_1\}$  on choices of period 1 as described by equation (II.2) is derived as follows. Individuals maximise their period one utility. By assumption, they only consider period 1 when choosing

period 1 pro-environmental effort. As such, they solve the following:

$$\partial_{x_1} U_1 - c_1 = 0 \quad (\text{A.1})$$

Where  $x_1^*$  is the solution to equation (A.1). It is a function of  $I_0, E_0, x_0, \eta_1, \epsilon_1,$  and  $c_1$ .

Using the implicit function theorem, we differentiate (A.1) with respect to  $\theta_1 \in \{\eta_1, \epsilon_1\}$ :

$$\partial_{x_1 x_1} U_1 \frac{\partial x_1^*}{\partial \theta_1} + \partial_{x_1 \theta_1} U_1 = 0 \Leftrightarrow \frac{\partial x_1^*}{\partial \theta_1} = \frac{\partial_{x_1 \theta_1} U_1}{-\partial_{x_1 x_1} U_1} > 0$$

The same reasoning applies when  $\theta_1 = c_1$ . ■

*Proof. (Spillovers of policies)* The optimal level of pro-environmental efforts at period 2 is a function of the choices of period one. Expressing the side effects of a policy at period one on choices of period two as in equation (II.3) amounts to using the implicit function theorem and differentiating the first order conditions of the maximisation programme at period 2 with respect to  $\theta_1 \in \{\eta_1, \epsilon_1\}$ :

$$\begin{aligned} \partial_{x_2 x_2} U_2 \frac{dx_2^*}{d\theta_1} + \partial_{x_2 x_1} U_2 \frac{\partial x_1^*}{\partial \theta_1} + \partial_{x_2 \theta_1} U_2 &= 0 \\ \Leftrightarrow \frac{dx_2^*}{d\theta_1} &= \frac{1}{-\partial_{x_2 x_2} U_2} \left[ \partial_{x_2 x_1} U_2 \frac{\partial x_1^*}{\partial \theta_1} + \partial_{x_2 \theta_1} U_2 \right] \end{aligned}$$

The same reasoning applies when  $\theta_1 = c_1$ . ■

*Proof. (Proposition 3)* First, model (III.4) can be rewritten in a reduced form as below:

$$x_{2i} = \bar{\alpha} + \bar{\beta}_1 \cdot c_{1i} + \bar{\beta}_2 \cdot \theta_{1i} + \bar{\epsilon}_i \quad (\text{A.2})$$

Where:

$$\bar{\alpha} = \alpha' + \beta^{BS} \cdot \alpha \quad \bar{\beta}_1 = \beta^{BS} \cdot \beta_1 \quad \bar{\beta}_2 = \beta^C + \beta^{BS} \cdot \beta_2 \quad (\text{A.3})$$

This implies that:

$$\beta^{BS} = \frac{\bar{\beta}_1}{\beta_1} \quad \beta^C = \bar{\beta}_2 - \frac{\bar{\beta}_1}{\beta_1} \beta_2 \quad (\text{A.4})$$

Using ordinary least square, we can show that the coefficients of model (A.2) are equal to:

$$\bar{\beta}_1 = \frac{\sigma_{2c}\sigma_\theta - \sigma_{2\theta}\sigma_{\theta c}}{\sigma_c\sigma_\theta - \sigma_{\theta c}^2} \quad \bar{\beta}_2 = \frac{\sigma_{2\theta}\sigma_c - \sigma_{2c}\sigma_{\theta c}}{\sigma_c\sigma_\theta - \sigma_{\theta c}^2}$$

Where:

$$\sigma_{2\theta} = \text{cov}(\mathbf{x}_2, \boldsymbol{\theta}_1) \quad \sigma_{\theta c} = \text{cov}(\boldsymbol{\theta}_1, \mathbf{c}_1)$$

And  $\sigma_\theta$  and  $\sigma_c$  denote respectively the variance of  $\boldsymbol{\theta}_1$  and  $\mathbf{c}_1$ . Furthermore, using ordinary least square, we can show that the coefficients of model (III.2) are equal to:

$$\beta^{ME} = \frac{\sigma_{1\theta}}{\sigma_\theta} \quad \beta^{SE} = \frac{\sigma_{2\theta}}{\sigma_\theta}$$

Similarly, the coefficients of the first stage of model (III.4) are equal to:

$$\beta_1 = \frac{\sigma_{1c}\sigma_\theta - \sigma_{1\theta}\sigma_{\theta c}}{\sigma_c\sigma_\theta - \sigma_{\theta c}^2} \quad \beta_2 = \frac{\sigma_{1\theta}\sigma_c - \sigma_{1c}\sigma_{\theta c}}{\sigma_c\sigma_\theta - \sigma_{\theta c}^2}$$

Injecting these expressions into expression  $\Xi = \beta^C + \beta^{BS} \cdot \beta^{ME} - \beta^{SE}$ , one can show that

$$\Xi = 0 \Leftrightarrow \beta^{SE} = \beta^C + \beta^{BS} \cdot \beta^{ME}$$

■

## B Appendix: Machine Learning Procedure

**Gradient tree boosting:** Let  $\{(x_i, y_i)\}_{i=1}^n$  be the training set with  $x_i$  the covariates of observation  $i$  and  $y_i$  its class. A decision tree is a function  $F$  which partitions the space of covariates into  $K$  regions  $\{R_1, \dots, R_K\}$ . It predicts a single class  $\hat{y}_k$  in each region, for  $k \in \{1, \dots, K\}$ :

$$F(x) = \sum_{k=1}^K \hat{y}_k \mathbf{1}_{R_k}(x)$$

Where  $\mathbf{1}_{R_k}(x)$  is the indicator function. We want to minimise  $L(y, F(x))$  where  $L$  is a loss function. This is done in  $M$  steps such that at each step  $m$ , we fit a function  $h_m \in \mathcal{H}$  to the "residuals" of the  $m - 1$  iteration such that:

$$F_m(x) = F_{m-1}(x) + v \cdot h_m(x, \delta_{km}) = F_{m-1}(x) + v \cdot \sum_{k=1}^K \delta_{km} a_{km} \mathbf{1}_{R_{km}}(x)$$

Where  $v$  is a shrinkage parameter reducing the speed at which the model is updated.  $a_{km}$  is the value predicted by  $h_m$  in the region  $R_{km}$ .  $h_m$  is called a base learner. The scalars  $\delta_{km}$  are set to minimise the loss function. For  $\gamma_{km} = \delta_{km} a_{km}$ :

$$\gamma_{km} = \arg \min_{\gamma} \sum_{x_i \in R_{km}} L(y_i, F_{m-1}(x_i) + \gamma)$$

The algorithm is defined as below:

**Algorithm:**

- Step 0: Choose a constant value  $\gamma$  such that:

$$F_0(x) = \arg \min_{\gamma \in \mathbb{R}} \left[ \sum_{i=1}^n L(y_i, \gamma) \right]$$

- Step  $m$ :

1. Compute the pseudo-residuals:

$$r_{im}(x_i) = -\frac{\partial L(y_i, F_{m-1}(x_1))}{\partial F_{m-1}(x_1)}$$

2. Fit a base learner  $h_m$  on the pseudo-residuals.

3. For each partition  $R_{km}$ , find the value  $\gamma_{km}$  such that:

$$\gamma_{km} = \arg \min_{\gamma} \sum_{x_i \in R_{km}} L(y_i, F_{m-1}(x_i) + \gamma)$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + v \cdot \sum_{k=1}^K \gamma_{km} \mathbf{1}_{R_{km}}(x)$$

- Step  $M$ : Output function  $F_M(x)$ .



**Tuning of hyperparameters:** The hyper-parameters we use in this paper are the following:

- The shrinkage parameter  $v$  is set to 0.01. Small values allow an improvement in performance by "forcing" the algorithm to learn slower.
- The bagging fraction is set to 0.5, meaning that 50% of the training observations are randomly drawn at each iteration to train the next tree expansion. Discarding

half of the observations reduces the over-fitting risk and improves computation speed.

- The minimal number of observations in each terminal node  $R_{km}$  is set to 50 when oversampling the *unwilling* and the *transitioned* and 10 in the case without oversampling. Splits leading to nodes with numbers of observations below this threshold are discarded. This parameter is tuned using grid search.
- The size of trees  $K$  is set to 7 when oversampling the *unwilling* and the *transitioned* and 8 in the case without oversampling. The higher this number, the more numerous the interactions between covariates (the "deeper" the tree). This parameter is tuned using grid search.
- The number of trees fitted  $M$  is set to 500 when oversampling and 450 in the case without oversampling. The lower the shrinkage parameter, the higher the number of trees has to be. This parameter is tuned using grid search.

In estimating the performances of GBM, we perform nested  $10 \times 10$  cross-validation. Namely, we randomly split the training set into ten subsets. First, the algorithm is fitted on nine subsets out of 10. Second, prediction errors are computed by comparing predictions made using the 10<sup>th</sup> subset data with respondents' actual answers. We repeat the first and the second steps ten times, each time with a new subset, to compute the prediction errors. This process is said to be nested as, at each step, the nine subsets used to fit the model are further split into ten subsets to tune the above hyperparameters. The process to select the hyperparameters maximising the prediction performances of the algorithm is similar to the process described at the beginning of this paragraph to estimate the algorithm's performance. Here, the performance met-



ric used is the average F1 score. Results are similar when using other metrics, such as Cohen’s Kappa.

**Performances estimation:** In total, we considered three different metrics to estimate the performances of GBM. First, for each type, we compute the share of individuals predicted to be of type  $i$  that are actually of type  $i$ . This measure is called *precision*. It tells us about the “purity” of our predicted classes. Precision should be higher than the share of respondents in type  $i$  over the total number of respondents to perform better than chance.<sup>26</sup> Here, GBM performs better than chance for each type and, on average, 1.9 times better than chance across all types (see Table 6).

Nevertheless, one can achieve high precision by excluding observations that are hard to predict. This is why we also look at *recall*, a measure of performance obtained by computing the proportion of individuals of type  $i$  correctly identified as being type  $i$ . This measure tells us about how “exhaustive” each predicted class is. With four types, a ratio above 25% indicates that the algorithm is performing better than chance.<sup>27</sup> The average recall rate of GBM is higher than 25% for the *unwilling type*, *hesitant type* and *trying type*, and slightly higher for the *transitioned type*. On average, GBM performs 1.6 times better than chance (see Table 6).

Ideally, we would like an algorithm yielding predicted types that are both “pure” and “exhaustive”. The F1 score is a measure encompassing these two aspects. It is the

---

<sup>26</sup>With four types, an algorithm doing as good as chance would produce a rate of true positives for type  $i$  to be  $\frac{n_i}{4}$ , where  $n_i$  is the number of individuals in type  $i$ . The rate of false positives would be  $\frac{n-n_i}{4}$  where  $n$  is the total number of respondents. Thus precision is equal to  $\frac{\frac{n_i}{4}}{\frac{n_i}{4} + \frac{n-n_i}{4}} = \frac{n_i}{n}$ .

<sup>27</sup>With an algorithm doing as good as chance, the rate of true positives is  $\frac{n_i}{4}$ , where  $n_i$  is the number of individuals in type  $i$ . The rate of false negatives is  $\frac{3n_i}{4}$ . Thus recall is equal to  $\frac{\frac{n_i}{4}}{\frac{n_i}{4} + \frac{3n_i}{4}} = \frac{1}{4}$ .

harmonic mean of precision and recall:

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

For our algorithm to perform better than chance, the average F1 score in each type should be higher than the thresholds displayed in Table 6.<sup>28</sup> Results in Table 6 confirm that GBM does better than chance for all types and on average 1.7 times better than chance across all types.

A closer look at Table 6 reveals that GBM over-classifies respondents as *hesitant* and under-classifies respondents as *unwilling* and *transitioned*. This explains the higher recall rate of the *hesitant* type and the higher precision rate of the *unwilling* and *transitioned* types. To correct this bias, we train another GBM algorithm where, this time, we over-sample the *unwilling* and *transitioned* types in the training set. Namely, we increase the sizes of these two sub-samples by drawing new observations with replacements from the original sub-samples. The new algorithm now seems to under-predict respondents to be *hesitant* in favour of the *unwilling* and *transitioned*. Although not statistically significant, over-sampling improves the overall recall rate of the model at the expense of precision and the F1 score. Furthermore, the relative sizes of each predicted type seem closer to these of the training set when over-sampling as measured by the Euclidian distance, although, here again, the difference is not statistically significant (see Table 6).

A last performance check consists of looking at whether the miss-classification errors of our two extreme types (*transitioned type* and *unwilling type*) occur in "adjacent" types.

---

<sup>28</sup>The minimum thresholds for the precision and the recall of an algorithm doing as good as chance are respectively  $\frac{n_i}{n}$ , and  $\frac{1}{4}$ . As such, the F1 score of this algorithm:  $2 \times \frac{\frac{n_i}{n} \times \frac{1}{4}}{\frac{n_i}{n} + \frac{1}{4}} = \frac{2 \times n_i}{n_i \times 4 + n}$ .

Table 6: Estimated performance of GBM

	Relative size of each type		Precision (in %)		Recall (in %)		F1 score (in %)	
	Training set	Predicted	Threshold	GBM	Threshold	GBM	Threshold	GBM
Unwilling	18.3	12 [24.2]	0.183	0.553 (0.03) [0.419*** (0.02)]	0.25	0.338 (0.02) [0.535*** (0.03)]	0.212	0.417 (0.02) [0.469 (0.02)]
Hesitant	39.4	53.5 [30.1]	0.394	0.482 (0.01) [0.491 (0.02)]	0.25	0.668 (0.02) [0.402*** (0.02)]	0.306	0.559 (0.01) [0.440*** (0.02)]
Trying	29.7	27.9 [27.0]	0.297	0.422 (0.02) [0.394 (0.01)]	0.25	0.397 (0.02) [0.346 (0.02)]	0.297	0.408 (0.02) [0.366 (0.01)]
Transitioned	12.6	6.6 [18.7]	0.126	0.469 (0.03) [0.352** (0.03)]	0.25	0.251 (0.03) [0.507*** (0.03)]	0.167	0.320 (0.03) [0.412** (0.03)]
	Euclidean distance (cross-validated)		Average					
	/	0.18 (0.02) [0.14 (0.01)]	0.25	0.481 (0.04) [0.414 (0.03)]	0.25	0.413 (0.03) [0.447 (0.03)]	0.246	0.426 (0.03) [0.421 (0.03)]

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: The columns labelled “threshold” contain the minimum performance threshold for each metric. Below these thresholds, GBM does worse than chance. Values in brackets correspond to the performance of GBM after over-sampling the *unwilling* and *transitioned* types. Stars indicate the results of simple t-tests to assess whether performances after re-sampling differ significantly from before.

Indeed, one would prefer to avoid using an algorithm that jumbles the *transitioned* and the *unwilling* types. Here, we estimate two sets of probabilities: the probability of being classified as type  $j$  whilst being of type  $i$ ,  $P(class = i | type = j)$ , and the probability of being of type  $j$  whilst being classified as type  $i$ ,  $P(type = i | class = j)$ . The first set of probabilities measures the model’s performance *ex-ante*: e.g., what is the probability that I will be classified in class  $k$  given my type? Symmetrically, the second set of probabilities gives us a measure of the model’s performance *ex-post*: e.g., what is the probability that I am of type  $k$  given how I was classified. The left panel of Figure 3 presents the estimated first set of probabilities, whilst the right panel presents the second. Overall, misclassification errors occur less often in non-adjacent categories. Furthermore, the left panel of Figure 3 indicates that over-sampling has increased the ability of the algorithm to correctly identify the *unwilling* and *transitioned* at the expense of its ability to identify the *hesitant* correctly. However, over-sampling has also slightly decreased its ability to produce pure predicted classes, as suggested by the right panel of Figure 3. In other words, over-sampling seems to make GBM better at detecting the

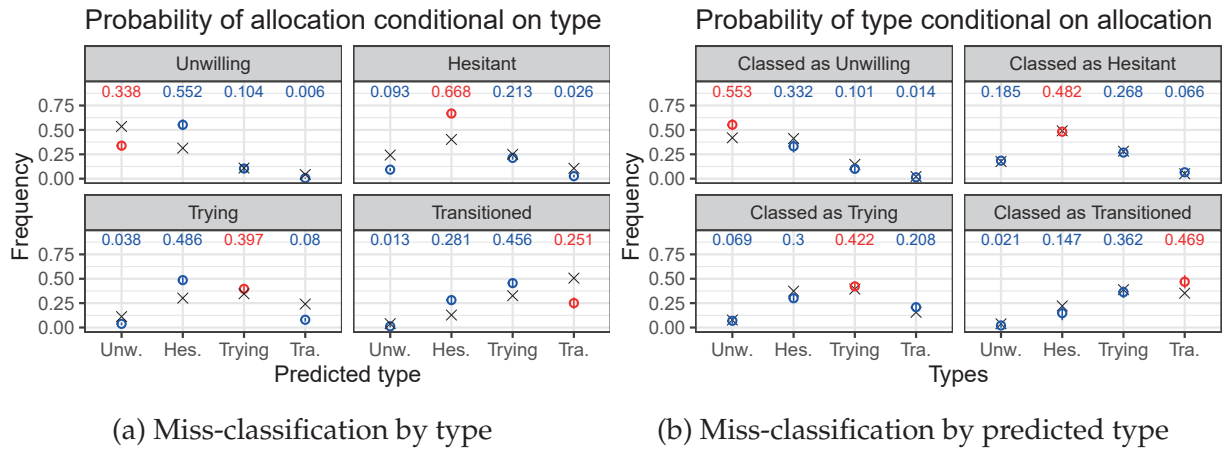


Figure 3: Frequency of miss-classification errors

Note: red dots correspond to the performance metric recall and precision on the right and left panels. Black crosses correspond to the estimates after re-sampling. 95% confidence intervals are represented by the vertical bars.

*unwilling* and *transitioned* types by simply increasing the number of respondents classified in these categories. We use the predictions obtained without over-sampling to carry out the main analysis.

**Predictive power of covariates:** The eighteen predictors used to train the GBM algorithm can be broadly grouped into four categories displayed in Table 7. First, sociological and economic characteristics of the area of residence of respondents account for 35.67% of the relative influence of the predictors. These variables are retrieved by merging information from the UK 2011 census data provided by the Office for National Statistics and the postcode respondents reported. Second, respondents' attitudes towards climate change and their knowledge of the environmental impact of food represent 33.87% of the relative influence of all the predictors. In this category, respondents' belief about whether acting against climate change is a moral duty has the greatest influence. In itself, it accounts for about 18% of the relative influence of the eighteen variables. Third, respondents' social-demographic characteristics represent 16.02% of the total influence of the predictors, followed by measures of respondents'

Table 7: Relative influence of each predictor

Category	Predictors	Relative influence (in %)
Social-demographics of residence area	Share of unemployed among actives in residence area	11.53
	Share of students in residence area	11.16
	Proportion of rural areas in residence area	7.00
	Share of UK/EU population in residence area	5.98
Belief and knowledge on the environment	Belief moral duty to act against climate change	17.93
	Knowledge of the carbon footprint of food	7.60
	Belief climate change is exaggerated	5.37
	Confidence in one's knowledge	2.97
Personal social-demographics	Age	8.76
	Education	3.90
	Sex	1.50
	Moved out of birth area	1.04
	Caucasian	0.43
Food preferences	Live in London	0.40
	Belief British food should be meat-based	4.63
	Online food ordering habits	3.95
	Preference for meat-based food	3.37
	Follows a specific diet	2.48

food preferences that account for 14.43% of this influence. We excluded two predictors that contained too many missing values: respondents' income and political beliefs. We include them back when testing for the robustness of our results. Readers interested in the influence of each predictor on the likelihood of being classified in one of the four types can refer to the partial dependency plots displayed in Figures 13, 14, and 15 in Appendix D.

## C Appendix: Robustness Checks

Table 8: Robustness checks of ATEs of the social norm message

Outcome	Chose vegetarian food (binary)					Food choice in kgCO <sub>2</sub> -eq		
	ITT with controls	CACE w/o controls	CACE with controls	Logit w/o controls	Logit with with controls	ITT with controls	CACE w/o controls	CACE with controls
Baseline	0.258*** (0.064)	0.136*** (0.012)	0.256*** (0.064)			22.874*** (4.295)	23.364*** (0.864)	22.889*** (4.297)
Social norm	0.057*** (0.016)	0.083*** (0.019)	0.070*** (0.020)	0.067*** (0.016)	0.056*** (0.016)	-2.211** (0.955)	-3.375** (1.137)	-2.698** (1.164)
Plant-intensive Chef's selection	0.124*** (0.016)	0.113*** (0.016)	0.123*** (0.016)	0.114*** (0.015)	0.121*** (0.016)	-7.697*** (0.961)	-7.802*** (0.929)	-7.639*** (0.961)
Num.Obs.	2454	2775	2454	2775	2454	2453	2775	2453
R2	0.113					0.081		

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Effect of the social norm message and of default menu allocation on food choices when controls are added and with non-linear probit estimation. We use the following lasso-selected controls to increase the precision of the estimates: level of hunger, how busy one is at the moment of taking the survey, knowledge of the environmental impact of food and confidence in one's knowledge, online food ordering frequency, belief that British food should be meat-based, preference for meat-based food, belief that climate change is exaggerated, belief that acting against climate change is a moral duty, income, sex, political orientation, education level and a dummy capturing the visual aspect of the menu. Robust standard errors are displayed in parentheses.

Table 9: Robustness checks of the behavioural spillovers and crowding-in/out effects of the social norm message I

Outcome	Decision to donate (binary)						
	Chose vegetarian food (binary)				Food choice in kgCO <sub>2</sub> -eq		
Food choice	2SLS with controls	Probit w/o controls	Probit with controls	MLE w/o controls	2SLS with controls	Probit w/o controls	Probit with controls
Baseline	0.110 (0.069)				0.292*** (0.075)		
Food choice	0.329** (0.156)	0.355** (0.163)	0.323** (0.153)	0.351*** (0.017)	-0.005** (0.002)	-0.005** (0.002)	-0.005** (0.002)
Social norm	-0.018 (0.021)	-0.016 (0.022)	-0.017 (0.021)	-0.016 (0.018)	-0.010 (0.020)	-0.006 (0.020)	-0.010 (0.020)
Num.Obs.	2603	2775	2603	2775	2603	2775	2603

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Behavioural spillovers and crowding-in/out effects of the social norm message on the decision to donate. We use the following lasso-selected controls to increase the precision of the estimates: belief that British food should be meat-based, preference for meat-based food, belief that climate change is exaggerated, belief that acting against climate change is a moral duty and political orientation. The second, third, sixth and seventh columns contain estimates obtained with a two-stage [Rivers and Vuong \(1988\)](#) probit estimation. The fourth column contains estimates obtained with maximum likelihood estimation *à la* [Evans and Schwab \(1995\)](#) with standard errors obtained using the delta method. The other standard errors are robust and displayed in parentheses.

Table 10: Robustness checks of the behavioural spillovers and crowding-in/out effects of the social norm message II

Outcome	Amount donated (in £)	
Food choice	Binary food choice	Food choice in kgCO2-eq
Specification	2SLS with controls	2SLS with controls
Baseline	-0.002 (0.570)	1.121* (0.615)
Food choice	2.058* (1.235)	-0.033* (0.020)
Social norm	-0.188 (0.172)	-0.136 (0.160)
Num.Obs.	2602	2602

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: Behavioural spillovers and crowding-in/out effects of the social norm message on the amount donated. We use the following lasso-selected controls to increase the precision of the estimates: belief that British food should be meat-based, preference for meat-based food, belief that climate change is exaggerated, belief that acting against climate change is a moral duty, age and political orientation. Robust standard errors are displayed in parentheses.

Table 11: Effect of food choices on perception of effort for the environment

Outcome	Perception of effort			
Food choice	Binary		In kgCO2-eq	
Specification	OLS w/o controls	OLS with controls	OLS w/o controls	OLS with controls
Baseline	2.973*** (0.020)	1.822*** (0.162)	3.153*** (0.021)	2.012*** (0.161)
Food choice	0.310*** (0.042)	0.269*** (0.046)	-0.006*** (0.001)	-0.005*** (0.001)
Num.Obs.	2775	2453	2775	2453
R2	0.020	0.111	0.026	0.116

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Note: Effect of food choices on the perception of having exerted an effort for the environment. We control for the default menus, exposure to the social norm message, the appearance of menus, self-reported level of hunger and hurry, knowledge of the environmental impact of food and confidence in one's knowledge, frequency of food online delivery, income, age, education, belief that British food should be meat-based, preference for meat-based food, belief that climate change is exaggerated, belief that acting against climate change is a moral duty, gender, and political orientation. Robust standard errors are displayed in parentheses.

Table 12: Test of Assumption 1

Outcome	Amount donated (in £)		Decision to donate (binary)		Amount donated (in £)		Decision to donate (binary)	
	Binary	In kgCO2-eq	Binary	In kgCO2-eq	Binary	In kgCO2-eq	Binary	In kgCO2-eq
Specification	OLS				2SLS			
Baseline	3.064*** (0.117)	3.728*** (0.138)	0.450*** (0.015)	0.522*** (0.017)	2.795*** (0.414)	4.128*** (0.651)	0.428*** (0.052)	0.555*** (0.081)
Food choice	1.274*** (0.285)	-0.022*** (0.004)	0.140*** (0.034)	-0.002*** (0.001)	2.678 (2.096)	-0.042 (0.033)	0.254 (0.260)	-0.004 (0.004)
Social norm	-0.195 (0.167)	-0.128 (0.190)	-0.019 (0.021)	0.008 (0.023)	-0.004 (0.607)	-0.377 (0.763)	-0.057 (0.077)	0.031 (0.095)
Food choice × Social norm	0.389 (0.377)	0.004 (0.006)	0.069 (0.045)	0.000 (0.001)	-0.711 (2.672)	0.015 (0.040)	0.186 (0.337)	-0.002 (0.005)
Num.Obs.	2775	2775	2775	2775	2775	2775	2775	2775
R2	0.025	0.015	0.023	0.015				

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Saturated model allowing to test Assumption 1. For the last four columns, variables capturing food choices are instrumented by the default menu allocation, and the variables corresponding to the interaction between food choices and the social norm nudge are instrumented by default menu allocation interacted with the social norm nudge. Robust standard errors are displayed in parentheses.

Table 13: Robustness checks of the ATEs of the social norm message conditional on predicted classes I

Specification	Nested OLS model						
	Chose vegetarian food (binary)						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	0.076*** (0.015)	0.079*** (0.022)	0.073*** (0.020)	0.056*** (0.018)	0.071*** (0.021)	0.067*** (0.018)	0.114*** (0.021)
Hesitant	0.103*** (0.024)	0.098*** (0.027)	0.098*** (0.025)	0.094*** (0.022)	0.073*** (0.024)	0.069*** (0.022)	0.038 (0.025)
Trying	0.116*** (0.025)	0.160*** (0.031)	0.162*** (0.030)	0.238*** (0.029)	0.215*** (0.031)	0.219*** (0.029)	0.125*** (0.032)
Transitioned	0.301*** (0.035)	0.421*** (0.067)	0.427*** (0.065)	0.492*** (0.080)	0.444*** (0.064)	0.451*** (0.058)	0.306*** (0.046)
Social norm × Unwilling	0.042* (0.023)	-0.018 (0.030)	0.016 (0.031)	0.013 (0.028)	-0.001 (0.030)	0.003 (0.026)	0.001 (0.030)
Social norm × Hesitant	q=0.070 0.027 (0.028)	q=0.554 0.037 (0.022)	q=0.599 0.029 (0.020)	q=0.654 0.056*** (0.020)	q=0.981 0.061*** (0.019)	q=0.911 0.073*** (0.020)	q=0.970 0.076*** (0.021)
Social norm × Trying	q=0.335 0.106*** (0.031)	q=0.102 0.095*** (0.034)	q=0.136 0.124*** (0.033)	q=0.005 0.076** (0.033)	q=0.002 0.069** (0.034)	q=0.001 0.061* (0.034)	q<0.001 0.088** (0.036)
Social norm × Transitioned	q=0.001 0.062 (0.043)	q=0.004 0.061 (0.081)	q<0.001 0.045 (0.080)	q=0.022 0.106 (0.102)	q=0.045 0.075 (0.079)	q=0.071 0.048 (0.072)	q=0.018 0.025 (0.056)
Num.Obs.	2730	2431	2730	2730	2730	2730	2730
R2	0.069	0.066	0.065	0.078	0.081	0.088	0.055

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: ATEs of the social norm message on the decision to choose vegetarian food. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).



Table 14: Robustness checks of the ATEs of the social norm message conditional on predicted classes II

Specification	Nested OLS model						
Outcome	Food choice in kgCO <sub>2</sub> -eq						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	25.004*** (1.515)	27.100*** (2.259)	25.206*** (2.134)	26.367*** (2.195)	25.826*** (2.201)	26.597*** (2.019)	24.119*** (1.801)
Hesitant	-2.675 (2.029)	-6.831*** (2.488)	-4.521* (2.340)	-4.858** (2.397)	-4.446* (2.404)	-5.368** (2.255)	-3.831* (2.059)
Trying	-9.046*** (1.908)	-12.064*** (2.538)	-8.822*** (2.447)	-12.683*** (2.446)	-11.082*** (2.475)	-11.470*** (2.315)	-5.995*** (2.278)
Transitioned	-12.592*** (2.034)	-19.403*** (2.910)	-16.847*** (2.922)	-16.701*** (3.525)	-16.733*** (3.068)	-18.594*** (2.680)	-13.532*** (2.396)
Social norm × Unwilling	-2.659 (2.114)	1.434 (3.370)	0.978 (3.165)	0.201 (3.204)	1.275 (3.233)	-0.814 (2.865)	-0.584 (2.521)
	q=0.208	q=0.672	q=0.755	q=0.950	q=0.700	q=0.772	q=0.823
Social norm × Hesitant	-3.822** (1.830)	-3.069** (1.408)	-3.189** (1.302)	-3.524*** (1.308)	-3.651*** (1.306)	-3.645*** (1.357)	-2.868** (1.373)
	q=0.035	q=0.029	q=0.014	q=0.008	q=0.005	q=0.007	q=0.039
Social norm × Trying	-0.535 (1.694)	-1.912 (1.622)	-2.448 (1.668)	-1.168 (1.500)	-1.382 (1.610)	-1.889 (1.604)	-4.907*** (1.863)
	q=0.753	q=0.241	q=0.151	q=0.428	q=0.390	q=0.238	q=0.010
Social norm × Transitioned	-2.844* (1.717)	0.157 (2.457)	0.090 (2.642)	-3.958 (3.288)	-2.087 (2.633)	0.719 (2.435)	0.886 (2.231)
	q=0.099	q=0.953	q=0.976	q=0.241	q=0.440	q=0.773	q=0.679
Num.Obs.	2730	2431	2730	2730	2730	2730	2730
R2	0.037	0.040	0.028	0.037	0.035	0.037	0.025

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: ATEs of the social norm message on the carbon footprint of respondents' food. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

Table 15: Robustness checks of the side effects of the social norm message conditional on predicted classes I

Specification	Nested OLS model						
Outcome	Amount donated (in £)						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	1.557*** (0.174)	1.477*** (0.256)	1.273*** (0.232)	1.327*** (0.245)	1.383*** (0.224)	1.161*** (0.224)	1.747*** (0.215)
Hesitant	1.458*** (0.256)	1.491*** (0.298)	1.649*** (0.270)	1.605*** (0.281)	1.605*** (0.266)	1.692*** (0.263)	1.435*** (0.263)
Trying	2.612*** (0.274)	3.109*** (0.338)	3.378*** (0.316)	3.259*** (0.330)	3.282*** (0.311)	3.761*** (0.308)	2.409*** (0.321)
Transitioned	3.340*** (0.326)	3.633*** (0.585)	3.924*** (0.581)	3.835*** (0.570)	2.858*** (0.507)	3.672*** (0.690)	2.946*** (0.415)
Social norm × Unwilling	0.638** (0.263)	0.287 (0.392)	0.693* (0.367)	0.316 (0.364)	0.519 (0.336)	0.811** (0.358)	0.313 (0.313)
	q=0.014	q=0.483	q=0.060	q=0.375	q=0.111	q=0.022	q=0.321
Social norm × Hesitant	0.110 (0.265)	0.181 (0.216)	0.161 (0.198)	0.132 (0.194)	0.167 (0.204)	0.181 (0.195)	0.052 (0.214)
	q=0.671	q=0.417	q=0.414	q=0.493	q=0.409	q=0.363	q=0.807
Social norm × Trying	-0.542* (0.302)	-0.528* (0.314)	-0.689** (0.303)	-0.510 (0.310)	-0.655** (0.308)	-0.838*** (0.295)	-0.192 (0.336)
	q=0.071	q=0.098	q=0.021	q=0.100	q=0.034	q=0.003	q=0.566
Social norm × Transitioned	-0.518 (0.371)	-0.385 (0.672)	-0.920 (0.681)	-0.625 (0.676)	-0.197 (0.592)	-0.410 (0.864)	-0.698 (0.470)
	q=0.166	q=0.565	q=0.174	q=0.359	q=0.741	q=0.636	q=0.140
Num.Obs.	2730	2431	2730	2730	2730	2730	2730
R2	0.062	0.053	0.053	0.055	0.051	0.062	0.038

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Total side effects of the social norm message on the amount donated. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

Table 16: Robustness checks of the side effects of the social norm message conditional on predicted classes II

Specification	Nested OLS model						
Outcome	Decision to donate (binary)						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	0.239*** (0.024)	0.219*** (0.034)	0.194*** (0.031)	0.199*** (0.032)	0.202*** (0.029)	0.180*** (0.030)	0.279*** (0.030)
Hesitant	0.223*** (0.034)	0.224*** (0.039)	0.248*** (0.036)	0.251*** (0.037)	0.257*** (0.035)	0.255*** (0.035)	0.192*** (0.035)
Trying	0.351*** (0.034)	0.408*** (0.042)	0.439*** (0.039)	0.411*** (0.041)	0.418*** (0.038)	0.478*** (0.039)	0.286*** (0.041)
Transitioned	0.409*** (0.039)	0.453*** (0.068)	0.473*** (0.066)	0.478*** (0.066)	0.388*** (0.062)	0.463*** (0.081)	0.350*** (0.050)
Social norm × Unwilling	0.084** (0.035)	0.026 (0.051)	0.087* (0.048)	0.046 (0.048)	0.075* (0.044)	0.121** (0.049)	0.035 (0.042)
Social norm × Hesitant	q=0.016 0.028 (0.035)	q=0.614 0.031 (0.028)	q=0.071 0.029 (0.026)	q=0.340 0.020 (0.026)	q=0.092 0.021 (0.027)	q=0.015 0.027 (0.026)	q=0.401 0.011 (0.027)
Social norm × Trying	q=0.433 -0.074** (0.036)	q=0.271 -0.057 (0.037)	q=0.257 -0.074** (0.036)	q=0.454 -0.038 (0.036)	q=0.435 -0.057 (0.036)	q=0.284 -0.091*** (0.034)	q=0.691 -0.009 (0.040)
Social norm × Transitioned	q=0.039 -0.022 (0.042)	q=0.120 0.012 (0.076)	q=0.045 -0.073 (0.076)	q=0.291 -0.055 (0.076)	q=0.103 0.003 (0.071)	q=0.007 0.011 (0.100)	q=0.822 -0.033 (0.055)
	q=0.604	q=0.876	q=0.325	q=0.481	q=0.972	q=0.910	q=0.553
Num.Obs.	2730	2431	2730	2730	2730	2730	2730
R2	0.066	0.058	0.052	0.052	0.052	0.058	0.037

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Total side effects of the social norm message on the decision to donate. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

Table 17: Robustness checks of the crowding-out/in effect of the social norm message conditional on predicted classes I

Specification	Nested 2SLS model						
Outcome	Amount donated (in £)						
Food choice	Chose vegetarian food (binary)						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	1.424*** (0.202)	1.303*** (0.275)	1.137*** (0.247)	1.054*** (0.238)	1.198*** (0.265)	1.264*** (0.241)	1.529*** (0.262)
Hesitant	1.280*** (0.285)	1.277*** (0.321)	1.466*** (0.297)	1.511*** (0.281)	1.472*** (0.293)	1.483*** (0.280)	1.362*** (0.267)
Trying	2.411*** (0.312)	2.759*** (0.390)	3.076*** (0.376)	3.302*** (0.420)	2.865*** (0.422)	2.895*** (0.417)	2.170*** (0.360)
Transitioned	2.817*** (0.515)	2.713*** (0.800)	3.128*** (0.805)	2.724*** (0.931)	3.020*** (0.807)	2.059*** (0.768)	2.360*** (0.577)
Social norm × Unwilling	0.565** (0.268)	0.327 (0.399)	0.663* (0.370)	0.787** (0.358)	0.318 (0.368)	0.513 (0.340)	0.310 (0.316)
Social norm × Hesitant	q=0.034 0.063 (0.267)	q=0.411 0.101 (0.221)	q=0.075 0.106 (0.201)	q=0.028 0.074 (0.210)	q=0.393 0.020 (0.211)	q=0.129 0.038 (0.227)	q=0.327 -0.093 (0.236)
Social norm × Trying	q=0.803 -0.726** (0.334)	q=0.647 -0.737** (0.334)	q=0.591 -0.921*** (0.341)	q=0.730 -0.983*** (0.309)	q=0.923 -0.636** (0.322)	q=0.866 -0.763** (0.316)	q=0.695 -0.359 (0.356)
Social norm × Transitioned	q=0.030 -0.625 (0.381)	q=0.027 -0.519 (0.697)	q=0.009 -1.003 (0.695)	q=0.001 -0.615 (0.866)	q=0.044 -0.762 (0.679)	q=0.014 -0.282 (0.586)	q=0.298 -0.746 (0.467)
Food choice	q=0.097 1.735 (1.317)	q=0.445 2.186* (1.241)	q=0.153 1.863 (1.268)	q=0.485 1.927 (1.253)	q=0.264 1.835 (1.270)	q=0.620 1.771 (1.284)	q=0.105 1.915 (1.299)
	q=0.011	q=0.001	q=0.005	q=0.005	q=0.005	q=0.008	q=0.008
Num.Obs.	2730	2431	2730	2730	2730	2730	2730

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Crowding-out/in effect of the social norm message on the amount donated. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

Table 18: Robustness checks of the crowding-out/in effect of the social norm message conditional on predicted classes II

Specification	Nested 2SLS model						
Outcome	Amount donated (in £)						
Food choice	Food choice in kgCO2-eq						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	2.170*** (0.499)	2.437*** (0.608)	1.954*** (0.523)	1.900*** (0.525)	2.016*** (0.534)	2.066*** (0.544)	2.408*** (0.500)
Hesitant	1.393*** (0.261)	1.249*** (0.332)	1.527*** (0.285)	1.556*** (0.278)	1.487*** (0.296)	1.467*** (0.285)	1.330*** (0.273)
Trying	2.390*** (0.321)	2.682*** (0.417)	3.140*** (0.356)	3.406*** (0.382)	2.964*** (0.386)	2.988*** (0.375)	2.245*** (0.338)
Transitioned	3.031*** (0.402)	2.945*** (0.695)	3.469*** (0.658)	3.204*** (0.761)	3.388*** (0.656)	2.380*** (0.616)	2.575*** (0.482)
Social norm × Unwilling	0.573** (0.272)	0.337 (0.407)	0.719* (0.375)	0.817** (0.368)	0.350 (0.382)	0.498 (0.344)	0.297 (0.320)
Social norm × Hesitant	q=0.034 0.016 (0.276)	q=0.396 0.072 (0.226)	q=0.053 0.074 (0.207)	q=0.026 0.082 (0.208)	q=0.353 0.034 (0.207)	q=0.145 0.074 (0.216)	q=0.348 -0.027 (0.222)
Social norm × Trying	q=0.953 -0.555* (0.301)	q=0.747 -0.596* (0.318)	q=0.709 -0.755** (0.305)	q=0.695 -0.870*** (0.295)	q=0.871 -0.547* (0.310)	q=0.729 -0.703** (0.309)	q=0.906 -0.326 (0.348)
Social norm × Transitioned	q=0.066 -0.587 (0.371)	q=0.062 -0.379 (0.661)	q=0.015 -0.917 (0.675)	q=0.003 -0.521 (0.873)	q=0.081 -0.681 (0.677)	q=0.023 -0.178 (0.592)	q=0.345 -0.674 (0.462)
Food choice	q=0.111 -0.025 (0.019)	q=0.576 -0.035* (0.020)	q=0.169 -0.027 (0.018)	q=0.559 -0.028 (0.018)	q=0.302 -0.027 (0.018)	q=0.768 -0.026 (0.019)	q=0.146 -0.027 (0.019)
	q=0.012	q=0.001	q=0.004	q=0.002	q=0.005	q=0.008	q=0.005
Num.Obs.	2730	2431	2730	2730	2730	2730	2730

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Crowding-out/in effect of the social norm message on the amount donated. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

Table 19: Robustness checks of the crowding-out/in effect of the social norm message conditional on predicted classes III

Specification	Nested 2SLS model						
Outcome	Decision to donate (binary)						
Food choice	Chose vegetarian food (binary)						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	0.218*** (0.027)	0.191*** (0.037)	0.174*** (0.033)	0.164*** (0.032)	0.179*** (0.034)	0.184*** (0.031)	0.246*** (0.036)
Hesitant	0.196*** (0.038)	0.191*** (0.043)	0.221*** (0.040)	0.227*** (0.038)	0.231*** (0.039)	0.238*** (0.036)	0.180*** (0.036)
Trying	0.321*** (0.039)	0.353*** (0.050)	0.394*** (0.048)	0.408*** (0.053)	0.352*** (0.053)	0.360*** (0.052)	0.249*** (0.046)
Transitioned	0.330*** (0.063)	0.309*** (0.097)	0.354*** (0.098)	0.318*** (0.114)	0.355*** (0.097)	0.268*** (0.096)	0.260*** (0.071)
Social norm × Unwilling	0.073** (0.036)	0.032 (0.052)	0.082* (0.049)	0.118** (0.049)	0.046 (0.049)	0.074* (0.044)	0.035 (0.043)
Social norm × Hesitant	q=0.042 0.021 (0.035)	q=0.547 0.018 (0.029)	q=0.100 0.021 (0.026)	q=0.018 0.011 (0.027)	q=0.346 0.003 (0.027)	q=0.099 0.001 (0.029)	q=0.421 -0.012 (0.030)
Social norm × Trying	q=0.548 -0.102** (0.041)	q=0.519 -0.090** (0.039)	q=0.404 -0.108*** (0.040)	q=0.678 -0.114*** (0.036)	q=0.918 -0.057 (0.038)	q=0.963 -0.074** (0.037)	q=0.701 -0.035 (0.042)
Social norm × Transitioned	q=0.012 -0.039 (0.044)	q=0.025 -0.009 (0.081)	q=0.008 -0.085 (0.081)	q=0.001 -0.020 (0.103)	q=0.128 -0.076 (0.077)	q=0.041 -0.010 (0.071)	q=0.410 -0.040 (0.056)
Food choice	q=0.379 0.263 (0.166)	q=0.900 0.343** (0.155)	q=0.285 0.278* (0.160)	q=0.821 0.295* (0.158)	q=0.311 0.277* (0.160)	q=0.877 0.268* (0.161)	q=0.476 0.295* (0.164)
	q=0.003	q=0.000	q=0.001	q=0.000	q=0.002	q=0.002	q=0.001
Num.Obs.	2730	2431	2730	2730	2730	2730	2730

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Crowding-out/in effect of the social norm message on the amount donated. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).


















Table 20: Robustness checks of the crowding-out/in effect of the social norm message conditional on predicted classes IV

Specification	Nested 2SLS model						
Outcome	Decision to donate (binary)						
Food choice	Food choice in kgCO2-eq						
Algorithm	GBM 1	GBM 2	Random Forest	Ordered logit	Multinomial	LDA	QDA
Unwilling (baseline)	0.332*** (0.063)	0.369*** (0.077)	0.296*** (0.067)	0.293*** (0.068)	0.303*** (0.068)	0.305*** (0.069)	0.381*** (0.064)
Hesitant	0.213*** (0.035)	0.186*** (0.045)	0.230*** (0.038)	0.234*** (0.038)	0.233*** (0.039)	0.236*** (0.037)	0.175*** (0.037)
Trying	0.317*** (0.041)	0.341*** (0.054)	0.404*** (0.045)	0.424*** (0.049)	0.367*** (0.049)	0.374*** (0.047)	0.260*** (0.043)
Transitioned	0.362*** (0.049)	0.346*** (0.083)	0.405*** (0.077)	0.391*** (0.090)	0.410*** (0.077)	0.316*** (0.076)	0.293*** (0.059)
Social norm × Unwilling	0.074** (0.036)	0.034 (0.054)	0.091* (0.050)	0.122** (0.050)	0.051 (0.051)	0.072 (0.045)	0.033 (0.044)
Social norm × Hesitant	q=0.040 0.014 (0.036)	q=0.537 0.014 (0.029)	q=0.067 0.017 (0.027)	q=0.012 0.012 (0.027)	q=0.302 0.005 (0.027)	q=0.107 0.007 (0.028)	q=0.451 -0.001 (0.028)
Social norm × Trying	q=0.701 -0.076** (0.036)	q=0.648 -0.068* (0.037)	q=0.535 -0.083** (0.036)	q=0.647 -0.096*** (0.034)	q=0.856 -0.043 (0.036)	q=0.812 -0.065* (0.036)	q=0.962 -0.030 (0.042)
Social norm × Transitioned	q=0.035 -0.033 (0.042)	q=0.067 0.013 (0.074)	q=0.018 -0.072 (0.076)	q=0.005 -0.006 (0.101)	q=0.225 -0.064 (0.076)	q=0.073 0.005 (0.071)	q=0.472 -0.029 (0.054)
Food choice	q=0.418 -0.004 (0.002)	q=0.867 -0.006** (0.003)	q=0.346 -0.004* (0.002)	q=0.944 -0.004* (0.002)	q=0.393 -0.004* (0.002)	q=0.944 -0.004* (0.002)	q=0.585 -0.004* (0.002)
	q=0.003	q=0.000	q=0.001	q=0.001	q=0.001	q=0.001	q=0.001
Num.Obs.	2730	2431	2730	2730	2730	2730	2730

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: Crowding-out/in effect of the social norm message on the amount donated. The first column displays the results yielded by the predictions of the GBM algorithm with over-sampling and the second when political beliefs and income are used as predictors. The other columns display the results yielded by the predictions of the other algorithms. Robust standard errors are displayed in parentheses. P-values from randomisation tests with 5,000 re-sampling are displayed last (q).

## D Appendix: Supplementary Tables and Figures

			
<b>Ploughman's lunch (V)</b> <small>(693 kcal)</small> Freshly baked in-house wholegrain bread served with hard-boiled eggs, grilled vegetable, onions, pickles and fresh salad.	£10 	<b>Oxford style sausage (V)</b> <small>(682 kcal)</small> Vegetarian sausage flavoured with pepper, clove, sage, and mace. Served with mashed potatoes and mushroom sauce.	£10 
<b>Ploughman's lunch</b> <small>(871 kcal)</small> Freshly baked in-house wholegrain bread served with hard-boiled eggs, cheddar and ham, onions, pickles and fresh salad.	£14 	<b>Oxford sausage</b> <small>(986 kcal)</small> Veal and Pork sausage flavoured with pepper, clove, sage and mace. Served with mashed potatoes and mushroom sauce.	£12 
<b>Pie and mash (chicken)</b> <small>(681 kcal)</small> Chicken filling baked in puff pastry served with mashed potatoes and parsley sauce.	£12 	<b>Fish and chips</b> <small>(664 kcal)</small> Battered fillet of cod served with chips and mushy peas.	£11 
<b>Pie and mash (lamb)</b> <small>(925 kcal)</small> Lamb filling baked in puff pastry served with mashed potatoes and parsley sauce.	£16 	<b>Gammon steak</b> <small>(947 kcal)</small> Smoked pork roast served with chips and mushy peas.	£14 
		<b>Sunday roast (Ve)</b> <small>(607 kcal)</small> Roasted nut loaf served with Yorkshire pudding, roast potatoes and vegetables.	£10 
		<b>Sunday roast</b> <small>(947 kcal)</small> Rib of beef served with Yorkshire pudding, roast potatoes and vegetables.	£18 
		<b>Shepherd's pie (Ve)</b> <small>(619 kcal)</small> Minced vegetable with mashed potato crust.	£10 
		<b>Shepherd's pie</b> <small>(925 kcal)</small> Minced lamb with mashed potato crust.	£16 
<i>All our dishes are homemade with local ingredients.</i>			
 = Completely climate-friendly  = Somewhat climate-friendly  = Slightly climate-friendly  = Not climate-friendly at all			
V = vegetarian   Ve=vegan			




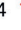













			
<b>Eggs and grilled vegetable (V)</b> <small>(693 kcal)</small> Freshly baked in-house wholegrain bread served with hard-boiled eggs, grilled vegetable, onions, pickles and fresh salad.	£10 	<b>Flavoured sausage (V)</b> <small>(682 kcal)</small> Vegetarian sausage flavoured with pepper, clove, sage, and mace. Served with mashed potatoes and mushroom sauce.	£10 
<b>Eggs, cheddar and ham</b> <small>(871 kcal)</small> Freshly baked in-house wholegrain bread served with hard-boiled eggs, cheddar and ham, onions, pickles and fresh salad.	£14 	<b>Flavoured sausage</b> <small>(986 kcal)</small> Veal and Pork sausage flavoured with pepper, clove, sage and mace. Served with mashed potatoes and mushroom sauce.	£12 
<b>Chicken pastry</b> <small>(681 kcal)</small> Chicken filling baked in puff pastry served with mashed potatoes and parsley sauce.	£12 	<b>Fillet of cod</b> <small>(664 kcal)</small> Fillet of cod in crispy coating served with chips and mashed peas.	£11 
<b>Lamb pastry</b> <small>(925 kcal)</small> Lamb filling baked in puff pastry served with mashed potatoes and parsley sauce.	£16 	<b>Smoked pork roast</b> <small>(947 kcal)</small> Smoked pork roast served with chips and mashed peas.	£14 
		<b>Roasted nut cake (Ve)</b> <small>(607 kcal)</small> Roasted nut cake served with popover, roast potatoes and vegetables.	£10 
		<b>Rib of beef</b> <small>(947 kcal)</small> Rib of beef served with popover, roast potatoes and vegetables.	£18 
		<b>Vegetable in potato crust (Ve)</b> <small>(619 kcal)</small> Minced vegetable with mashed potato crust.	£10 
		<b>Lamb in potato crust</b> <small>(925 kcal)</small> Minced lamb with mashed potato crust.	£16 
<i>All our dishes are homemade with local ingredients.</i>			
 = Completely climate-friendly  = Somewhat climate-friendly  = Slightly climate-friendly  = Not climate-friendly at all			
V = vegetarian   Ve=vegan			

Figure 4: Full menus

Note: two versions of the menus shown to participants. In total, we had 24 versions of the full menu in which we varied the ordering (12 versions) of the items and the menu's appearance (2 versions).



Table 21: Characteristics of the food items

Dish name	Main ingredients	Carbon footprint	Label colour
Eggs and grilled vegetable Ploughman's lunch	Eggs and vegetables	3.25	Green
Flavoured sausage Oxford style sausage	Beans	0.80	Green
Vegetable in potato crust Shepherd's pie	Vegetables	1.60	Green
Roasted nut cake Sunday roast	Nuts	2.00	Green
Chicken pastry Pie and mash	Chicken	5.40	Yellow
Fillet of cod Fish and chips	Fish	5.40	Yellow
Smoked pork roast Gammon steak	Pork	7.90	Orange
Eggs, cheddar and ham Ploughman's lunch	Ham and cheese	23.88	Red
Flavoured sausage Oxford sausage	Veal and pork	38.35	Red
Lamb in potato crust Pie and mash	Lamb	64.20	Red
Lamb pastry Shepherd's pie	Lamb	64.20	Red
Rib of beef Sunday roast	Beef	68.80	Red

Note: based on its carbon intensity, each dish is categorised in one of four categories, corresponding to the carbon footprint labels. Carbon footprints are computed based on the main ingredients of the dishes, using [Scarborough et al. \(2014\)](#)'s estimates. When dishes have more than one ingredient, we take the average between the two.

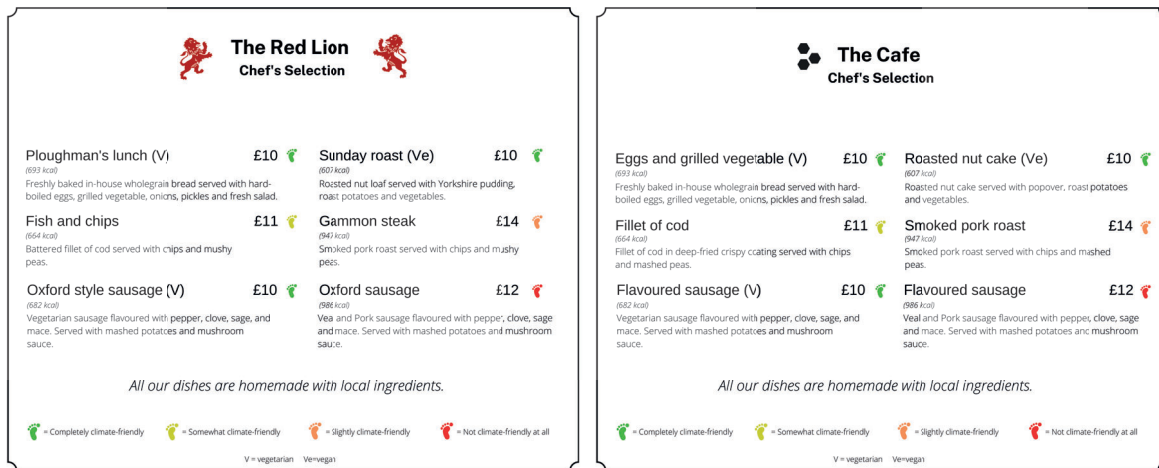


Figure 5: Plant-intensive default menus

Note: these are the two versions of the plant-intensive default menus shown to participants.

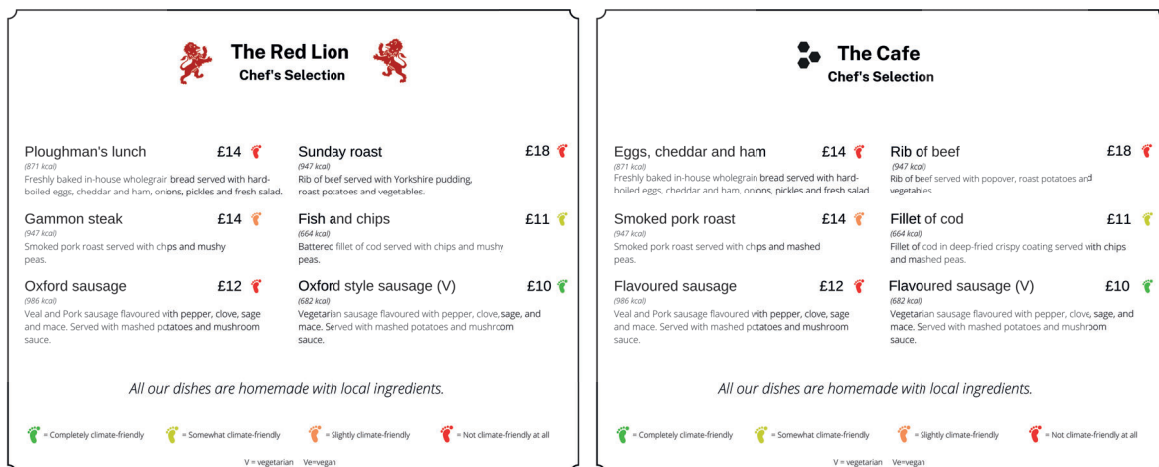


Figure 6: Meat-intensive default menus

Note: These are the two versions of the meat-intensive default menus shown to participants.

Table 22: Descriptive statistics

	Control group (n=1384)	Social norm group (n=1391)	p-value
<b>Age</b>			0.139
Mean	38.6 years old	37.9 years old	
Median	36 years old	35 years old	
<b>Income</b>			0.920
< £10,000	18.6%	17.9%	
£10,000 - £15,999	11.5%	12.5%	
£16,000 - £19,999	11.3%	10.8%	
£20,000 - £29,999	27.2%	28.1%	
£30,000 - £39,999	16.2%	14.9%	
£40,000 - £49,999	8.5%	8.4%	
£50,000 - £69,999	4.5%	4.6%	
£70,000 - £89,999	1.5%	1.9%	
£90,000 - £119,999	0.6%	0.5%	
£120,000 - £149,999	0.2%	0.2%	
More than £150,000	0.0%	0.2%	
<b>Gender</b>			0.450
Female	48.3%	51.0%	
Male	50.7%	48.2%	
Other	1.0%	0.7%	
<b>Education</b>			0.961
No education	0.1%	0.1%	
Primary education	0.2%	0.1%	
Lower secondary education	2.5%	2.6%	
Upper secondary education	22.6%	21.9%	
Post-secondary non-tertiary education	15.6%	15.0%	
Short-cycle tertiary education	5.5%	6.6%	
Bachelor or equivalent	40.2%	39.4%	
Master or equivalent	11.9%	12.9%	
Doctoral or equivalent	1.5%	1.4%	

Note: descriptive statistics per treatment group. A Wilcoxon test is carried out to test the difference in age between the treatment and the control group. A Chi-square test is used to test the difference in gender. Trend tests are carried out to test the differences in education and income between the two groups.

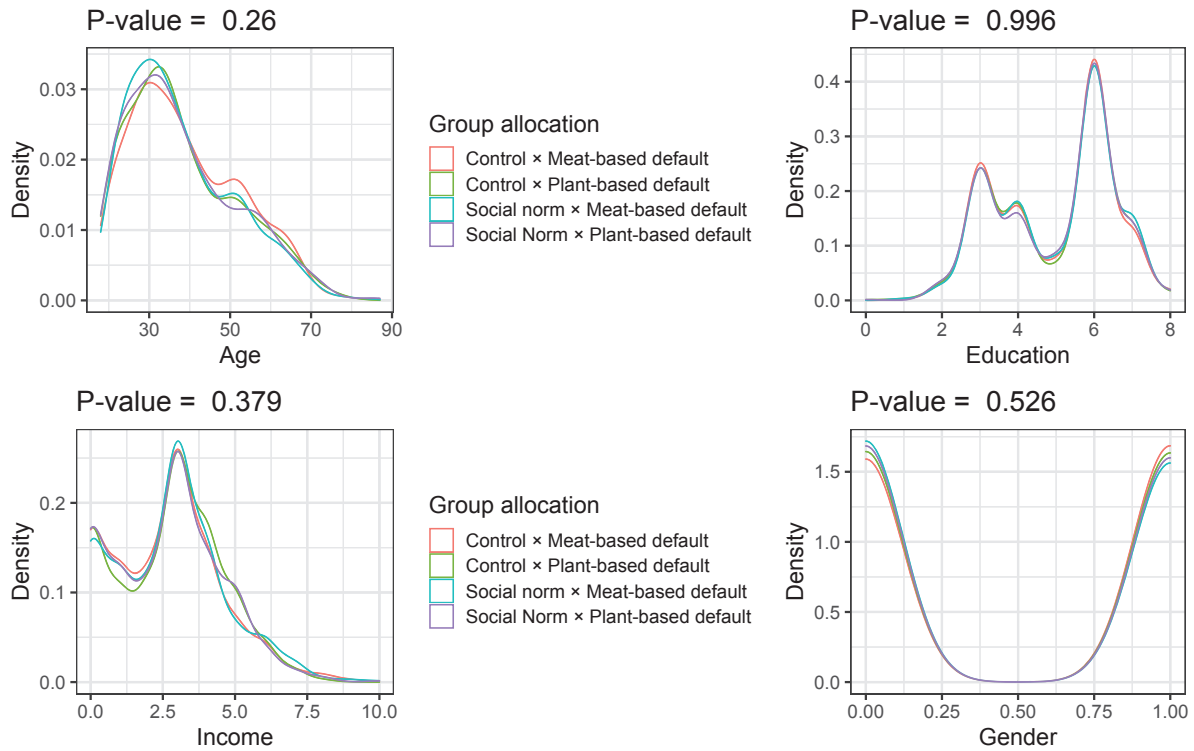


Figure 7: Distribution of the main covariates by treatment group

Note: density plots of age, education, income and gender across the four treatment groups of the *inference sample*. For education, 0 means "No education", and 8 means "PhD or equivalent". For Income, 0 means "less than £10k" and 10 means "more than £150k". For gender, 0 means female, and 1 means male.

Table 23: Descriptive statistics

<b>Main covariates</b>	
<b>Age</b>	
Mean	38 years old
Min	18 years old
Max	87 years old
SD	13.59 years old
<b>Income</b>	
Missing	339
< £10,000	969 (18.6%)
£10,000 - £15,999	673 (12.9%)
£16,000 - £19,999	580 (11.1%)
£20,000 - £29,999	1446 (27.7%)
£30,000 - £39,999	793 (15.2%)
£40,000 - £49,999	405 (7.8%)
£50,000 - £69,999	224 (4.3%)
£70,000 - £89,999	77 (1.5%)
£90,000 - £119,999	33 (0.6%)
£120,000 - £149,999	12 (0.2%)
More than £150,000	6 (0.1%)
<b>Gender</b>	
Missing	1
Female	2771 (49.9%)
Male	2736 (49.2%)
Agender	1 (0.0%)
Non-binary / third gender	42 (0.8%)
Trans woman	1 (0.0%)
Prefer not to say	5 (0.1%)
<b>Education</b>	
Missing	29
No education	2 (0.0%)
Primary education	12 (0.2%)
Lower secondary education	137 (2.5%)
Upper secondary education	1287 (23.3%)
Post-secondary non-tertiary education	853 (15.4%)
Short-cycle tertiary education	321 (5.8%)
Bachelor or equivalent	2166 (39.2%)
Master or equivalent	663 (12.0%)
Doctoral or equivalent	87 (1.6%)

Note: distribution of the main covariates across the 5,557 participants to the experiment.

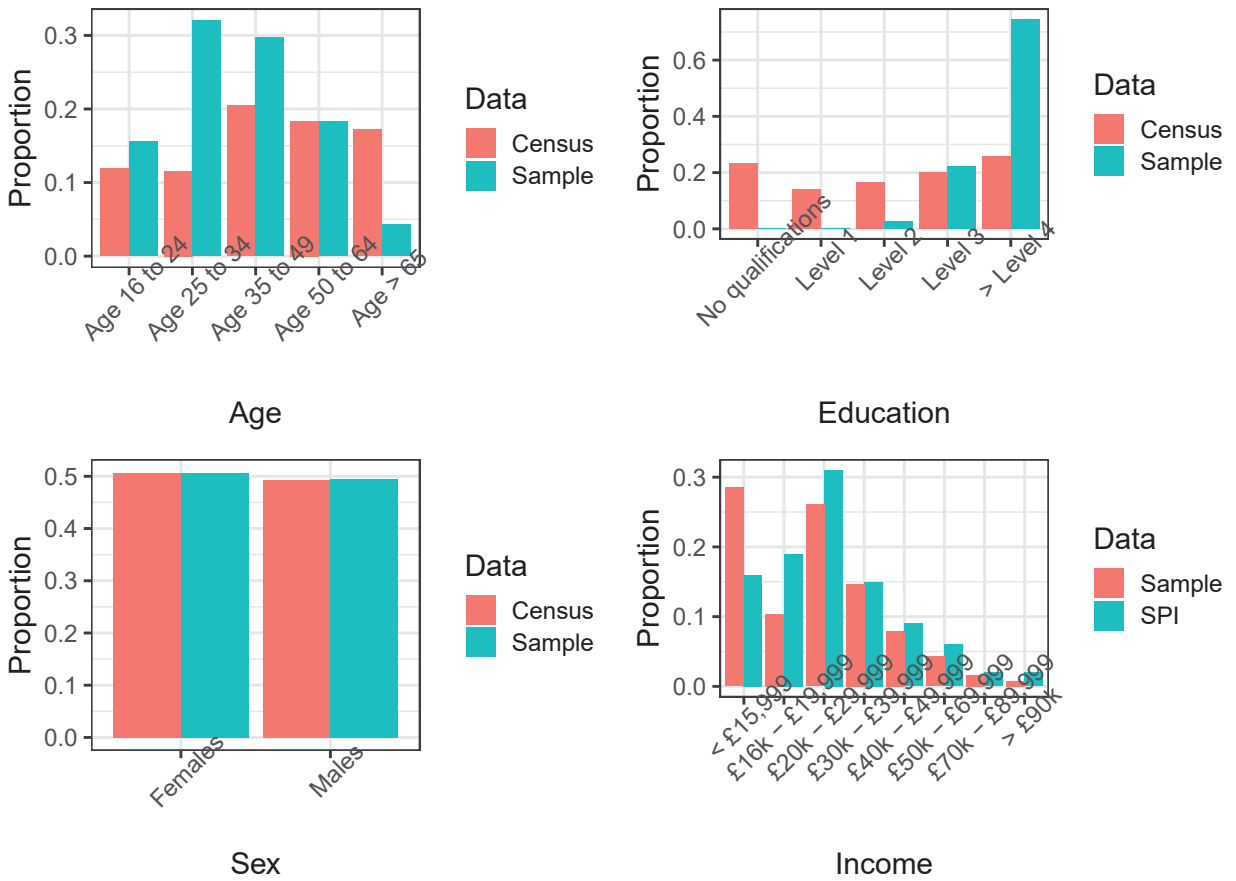


Figure 8: Comparison with UK population

Note: comparison of the distributions of the main covariates in the sample and the UK population. The data from the 2011 census is used to plot the distribution of age, sex and Education in the UK population. The 2020/2021 survey of personal income is used to plot the distribution of Income in the UK population.

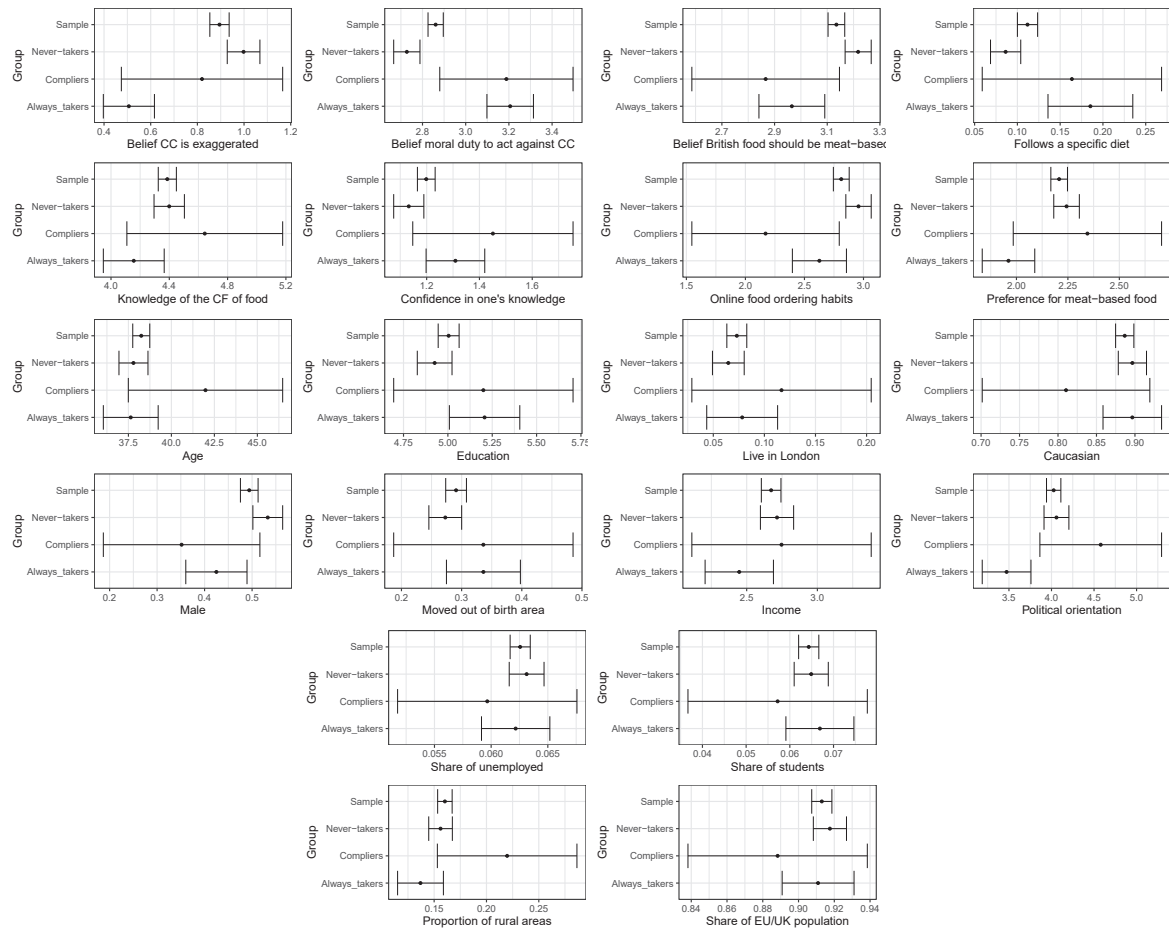


Figure 9: Profile of compliers

Note: we represent how the profile of compliers (those choosing vegetarian food when prompted to do so by the default nudge) differ from the rest of the sample, following [Marbach and Hangartner \(2020\)](#).

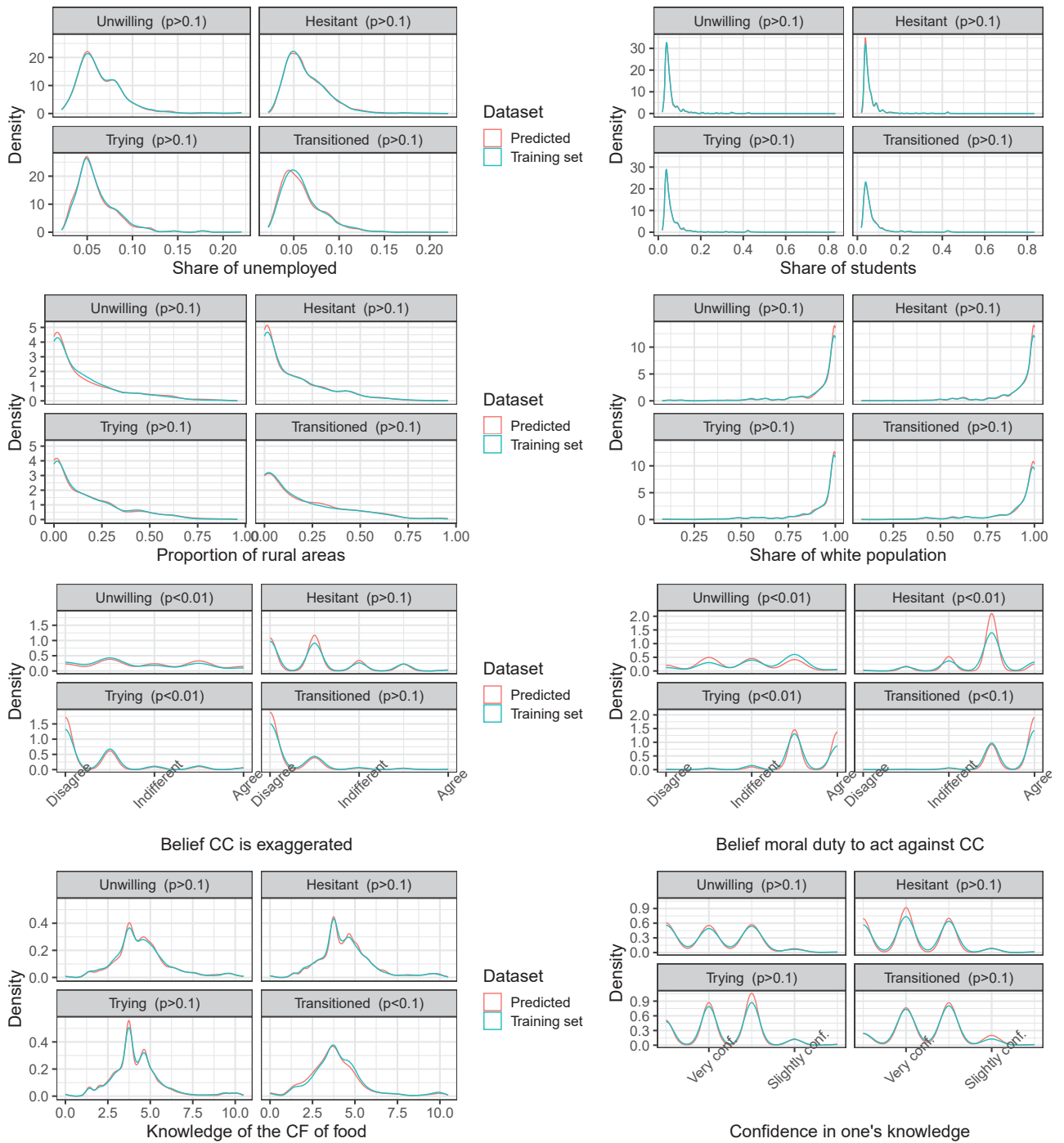


Figure 10: Distribution of the predictors by type I



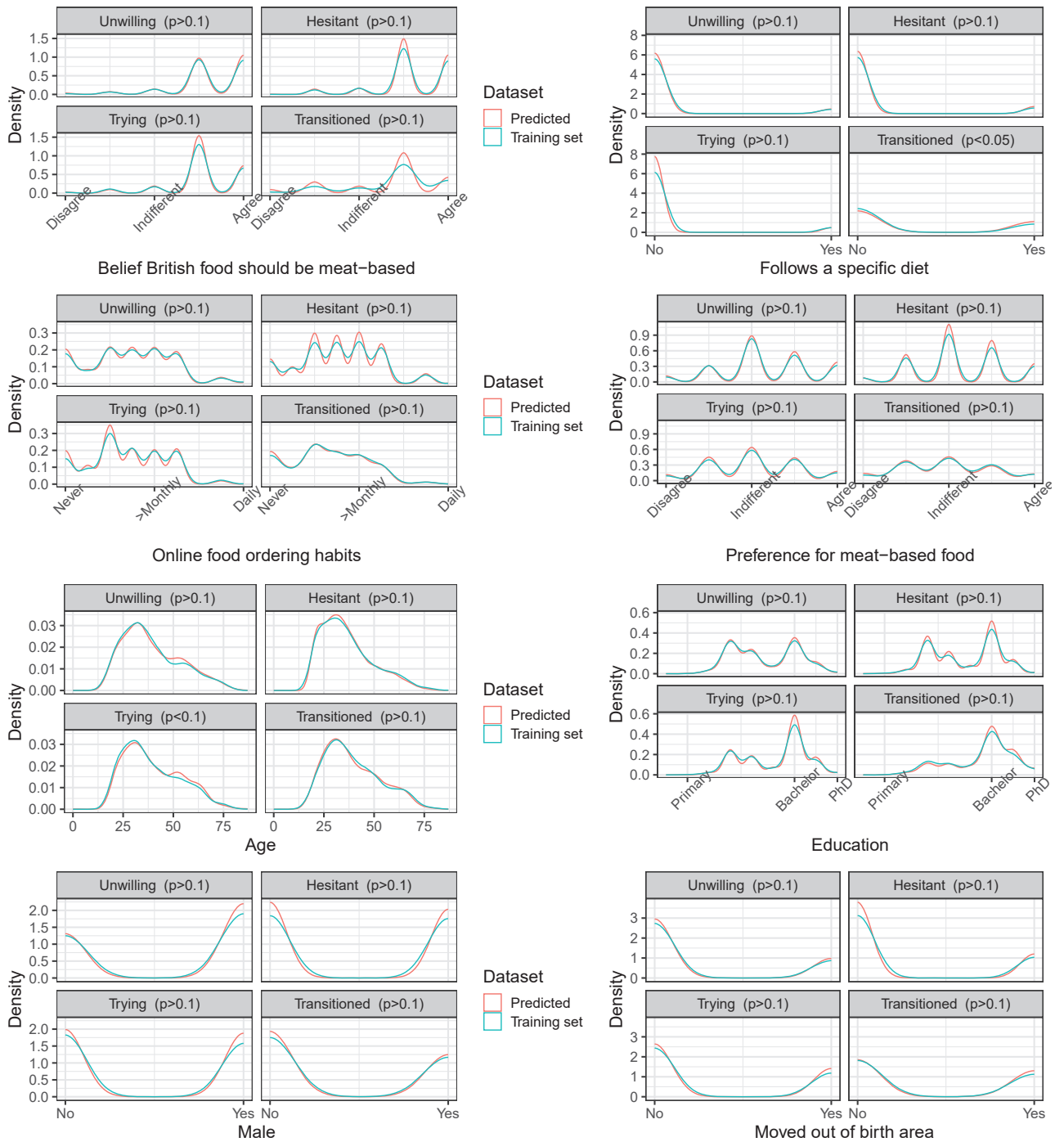


Figure 11: Distribution of the predictors by type II

Table 24: Profile of each predicted type

Covariates	Unwilling	Hesitant	Trying	Transitioned
Share of EU/UK population	0.004	0.003	-0.006	-0.004
Share of unemployed	0.003	0.004***	-0.006***	-0.006*
Proportion of rural areas	0.001	-0.023**	0.012	0.062**
Share of students	-0.001	-0.009**	0.009*	0.004
Belief moral duty to act against CC	-1.942***	-0.223***	0.957***	1.025***
Belief CC is exaggerated	1.598***	0.097	-0.696***	-0.794***
Knowledge of the CF of food	0.33**	0.442***	-0.43***	-1.008***
Confident in one's knowledge	-0.12	-0.265***	0.244***	0.484***
Age	2.851**	-4.787***	4.008***	0.601
Educated	-0.331**	-0.396***	0.369***	1.01***
Male	0.191***	-0.059*	0.023	-0.126**
Moved out of birth area	-0.034	-0.121***	0.112***	0.194***
Caucasian	-0.015	-0.003	0.021	-0.023
Live in London	-0.047***	-0.019	0.047***	0.012
Income	0.074	-0.332***	0.326***	0.06
Conservative	1.248***	0.177	-0.61***	-0.878***
Belief British food should be meat-based	0.174 **	0.157 ***	-0.036	-0.78 ***
Preference for meat-based food	0.224 **	0.192 ***	-0.229 ***	-0.456 ***
Follows a specific diet	-0.047 *	-0.005	-0.092 ***	0.397 ***
Order food online frequently	-0.089	0.541 ***	-0.441 ***	-0.549 ***

\* p < 0.1, \*\* p < 0.05, \*\*\* p < 0.01

Note: regression coefficients from linear models where each covariate is regressed on a dummy equal to 1 if respondents are classed in a given type, and zero otherwise. Coefficients, therefore, capture how different a given type is compared to the average of the sample. P-values are adjusted using Holmes-Bonferroni correction.

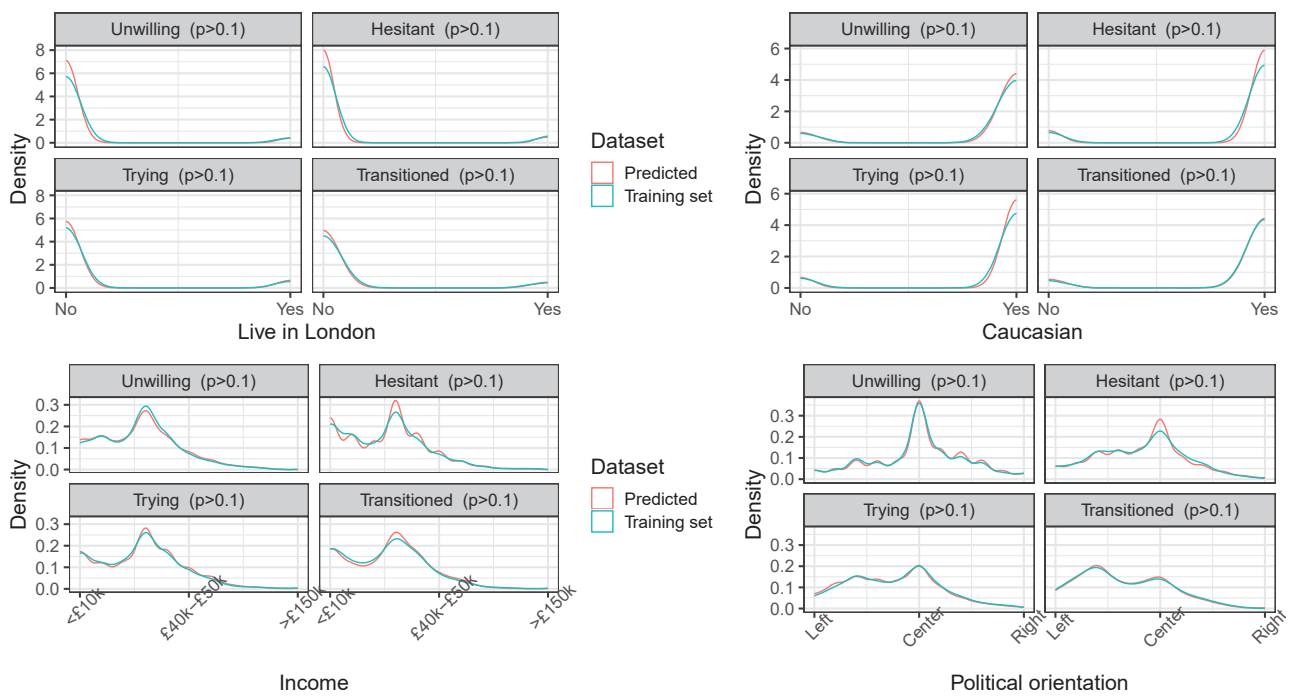


Figure 12: Distribution of the predictors by type III

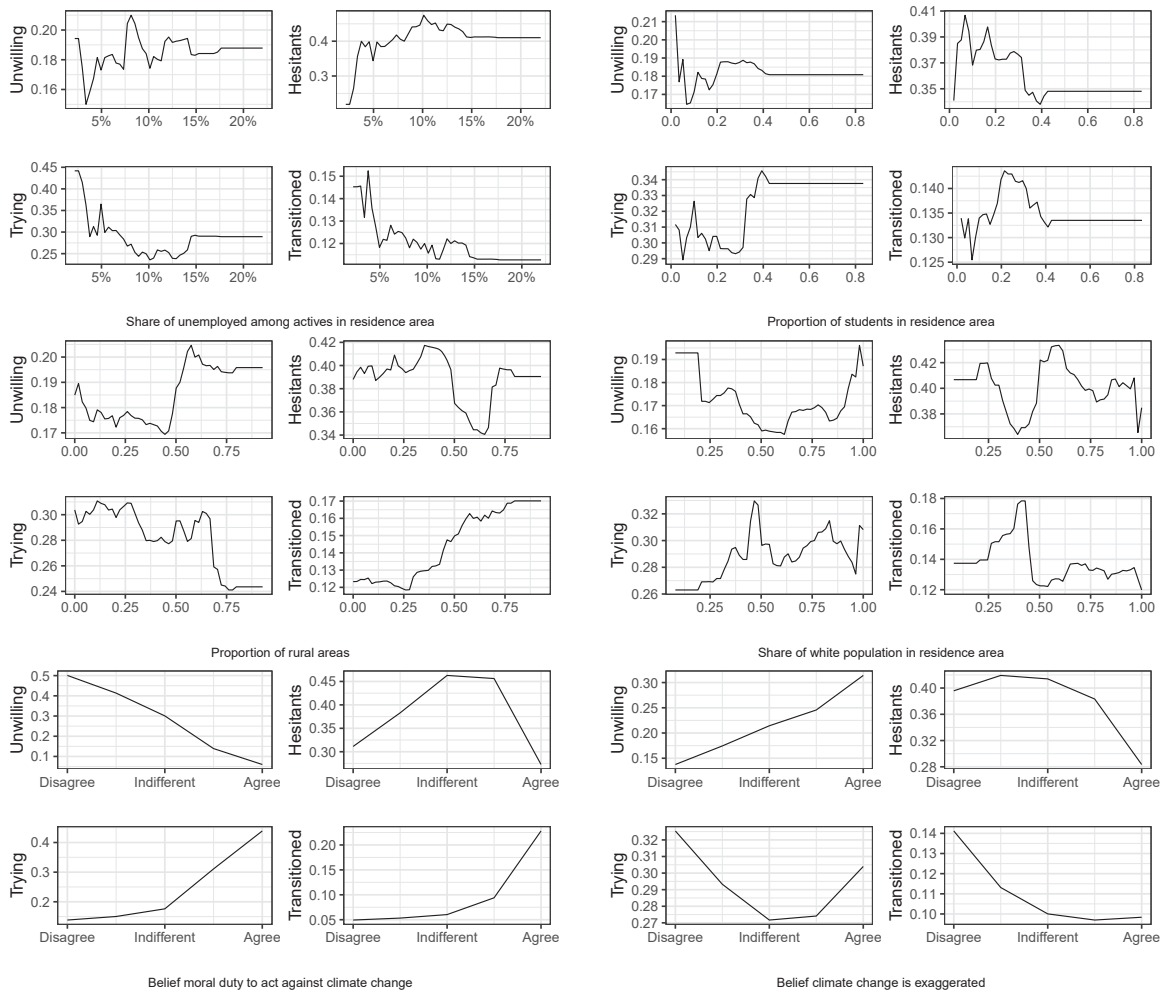


Figure 13: Partial dependence plots of the GBM algorithm I

Note: partial independence plots express visually the likelihood to be allocated to a given class against the values taken by a variable.

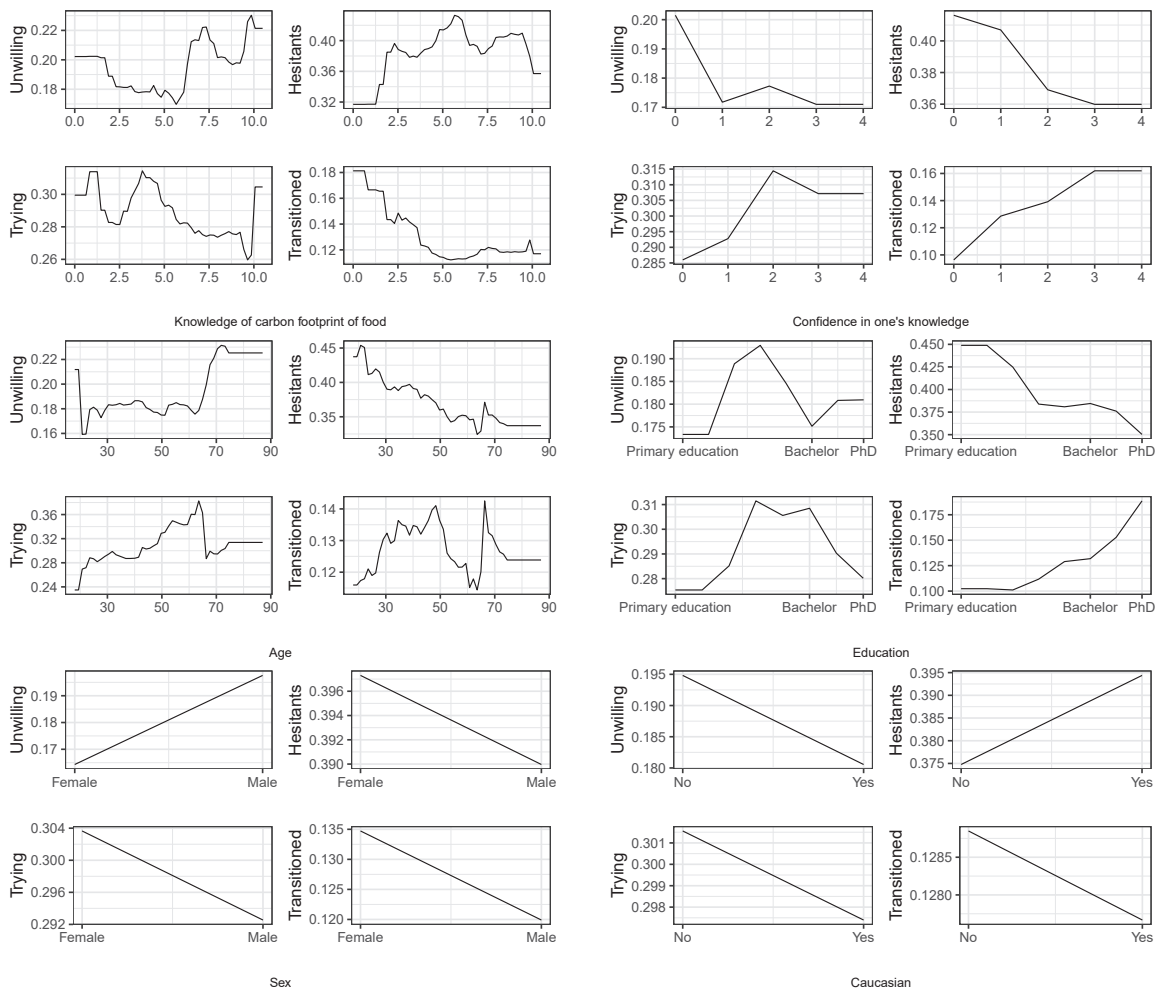


Figure 14: Partial dependence plots of the GBM algorithm II

Note: partial independence plots express visually the likelihood to be allocated to a given class against the values taken by a variable.

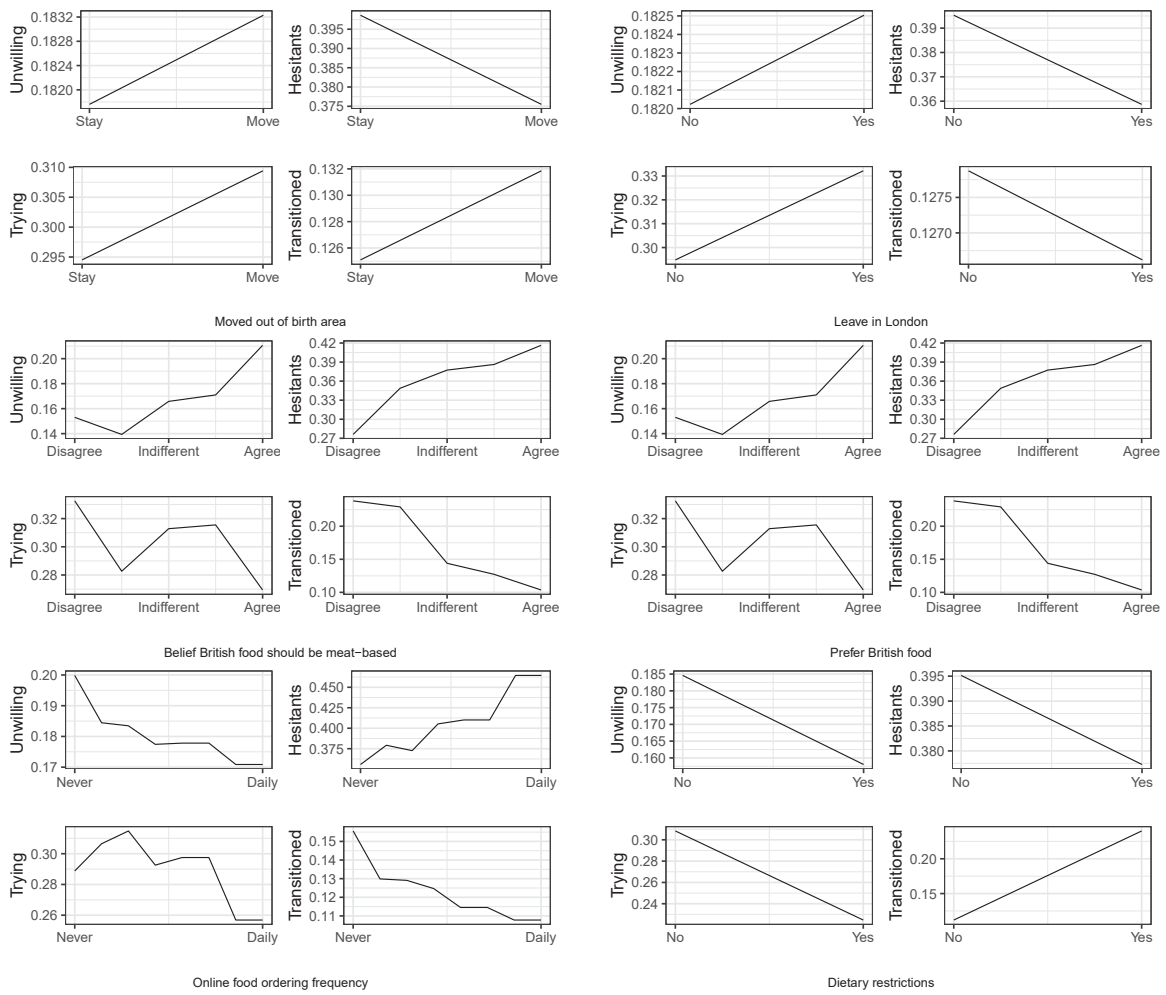


Figure 15: Partial dependence plots of the GBM algorithm III

Note: partial independence plots express visually the likelihood to be allocated to a given class against the values taken by a variable.