

# How Differently Do Farms Respond to Agri-environmental Policies? A Probabilistic Machine-Learning Approach

*Silvia Coderoni*

Department of Bioscience and Agro-Food and Environmental Technology, University of Teramo,  
Via Balzarini 1, 64100, Teramo (TE), Italy

email: [scoderoni@unite.it](mailto:scoderoni@unite.it)

*Roberto Esposti*

Department of Economics and Social Sciences, Università Politecnica delle Marche, Piazzale  
Martelli 8, 60121, Ancona (AN), Italy

email: [r.esposti@univpm.it](mailto:r.esposti@univpm.it)

*Alessandro Varacca*

Department of Economics and Social Sciences - DISES, Università Cattolica del Sacro Cuore, Via  
emilia Parmense 84, 29122, Piacenza (PC), Italy  
(*corresponding author*)

email: [varacca.alessandro@unicatt.it](mailto:varacca.alessandro@unicatt.it)

## Abstract

*This study evaluates to what extent farmers respond heterogeneously to the agri-environmental policies implemented within the European Common Agricultural Policy (CAP). Our identification and estimation strategy combines a theory-driven research design formalizing all possible sources of heterogeneity with a Bayesian additive regression trees algorithm. Results from a 2015-2018 panel of Italian farms show that the responsiveness to these policies may differ substantially across farms and farm groups. This suggests room for improvement in the implementation of these policies. We also argue that the specific features of the CAP call for a careful implementation of these empirical techniques.*

**Keywords:** Common Agricultural Policy, Agri-Environmental Measures, Heterogeneous Treatment Effects, Causal Machine Learning, Random Forests, Bayesian Additive Regression Trees.

**JEL Classification:** C21, Q15, Q51

# How differently do farms respond to agri-environmental policies? A probabilistic machine learning approach

## 1. Introduction

The Common Agricultural Policy (CAP) represents the primary ordinary policy instrument of the European Union (EU), at least in terms of budget share. Starting with the 1992 MacSharry reform, environmental and ecological concerns have increasingly become one of the major justifications for maintaining the CAP expenditure. Indeed, environmental policy objectives are likely to be the most relevant for European agriculture in the coming decades (Coderoni et al., 2021). Given the growing concerns about environmental and ecological issues and the resulting policy orientations, researchers are left to wonder how much farmers' behaviour has changed in response to the new greener CAP and what those responses are (Brown et al., 2021). Answering these questions is rather challenging, mainly because there is no univocal answer for the very large heterogeneity typically encountered in agriculture.

Since EU farmers are known for their distinctive diversity (Esposti, 2022a), we would typically expect equally diverse responses to these political shocks. Under this hypothesis, both academics and EU stakeholders have long advocated for a more targeted and tailored design of the EU policies (particularly CAP reforms – see Erjavec and Erjavec, 2015; Ehlers et al., 2021). However, such a task is challenging without a deeper understanding of whether and to what extent the potential recipients of such measures respond differently. As most parametric/semi-parametric (econometric) approaches to ex-post policy evaluation can only produce aggregate (i.e., average) responses or represent limited and pre-specified heterogeneity (see, for example, Esposti, 2017a; 2017b; Bertoni et al., 2020; Bartolini et al., 2021), the understanding of such heterogeneity has been so far rather limited.

Recent improvements in this field involve the use of specific causal inference (CI) methods (Imbens and Rubins, 2015) for framing the evaluation of a policy as a treatment effect discovery problem, which exploits counterfactual thinking to define the estimands of interest (Uehleke et al., 2022). Within the rapidly evolving literature, causal machine learning (CML) has started to gain attention as a useful extension to the more general CI framework, particularly when the objective of the evaluation regards highly complex and potentially heterogeneous responses to the treatment (Storm et al., 2020; Stetter et al., 2022). ML methods can be particularly beneficial when working with large heterogeneous samples characterised by many interacting variables and nonlinear relationships, but require suitable identifying assumptions and targeted technical adjustments (Chernozhukov et al., 2018; Hahn et al., 2018; Athey and Imbens, 2019). This means that off-the-shelf ML algorithms (i.e., common ML methods designed for predictive purposes) may, at best, represent one of several components in the CML toolbox. Given these premises, CML represents a suitable instrument for understanding how and to what extent the impact of recent CAP environmental policies varies across diverse farms. Our work fits within the very recent and fast developing empirical literature that deals with this issue. In particular, we aim at disentangling the causal effect of two alternative treatment options expressing two different implementations of the agri-environmental policy (AEP) within the 2014-2020 CAP reform. On the one hand, we consider farms that not only fulfil the basic eligibility conditions to benefit from the whole Direct Payment (DP), but also apply for Pillar 2 Agri-Environmental Measures (AEMs).<sup>1</sup> On the other hand, we consider farms that choose not to comply with the conditional requirements (i.e.: the so-called conditionality – see Section 2), thereby giving up the DP, and not take up any Pillar 2 AEM. We assume that the two treatments share the same control group, which consists of farms that only comply with the necessary environmental requirements to access the DP.

We begin by providing a theoretical background linking the determinants of AEP adoption by heterogeneous farmers to their production response and then, linking it to the potential outcomes framework. We exploit these conceptual underpinnings to define the relevant confounding variables

and treatments while providing a solid background for the necessary assumptions that characterise our identification strategy. The latter is grounded in the classical hypotheses that support most CI problems, included the Stable Unit Treatment Values Assumption (SUTVA) that may be problematic given the multiple-treatment nature of the AEMs. These hypotheses are coupled with flexible surface estimation via a CML algorithm known as Bayesian causal forests (BCF – Kapelner and Bleich, 2016; Carnegie et al., 2019; Hahn et al., 2020). Given their probabilistic nature, BCF can produce approximate posterior distributions for estimated heterogeneous treatment effects (HTE), allowing the introduction of uncertainty into group comparisons or, more generally, when transforming individual-level estimands. This feature represents a further original contribution of the present paper, as it may provide a useful improvement over other comparable ML methods for which inference is less straightforward (Stetter et al., 2022).

Our research is closely related with the recent analysis presented by Stetter et al. (2022), as both studies share a common objective of assessing the heterogeneous response of farmers to AEPs through CML techniques. Nonetheless, as elaborated in the preceding paragraphs and thoroughly discussed throughout the paper, our approach diverges from and extends upon their work in several fundamental aspects. These aspects encompass a more comprehensive delineation of the treatment set, a broader conceptualization of farmers' potentially heterogeneous response to AEPs, a distinct and relatively wider geographical coverage, and an investigation of the inherent limitations of conventional identification strategies employed in cross-sectional observational studies.

The remainder of the paper is organised as follows. Section 2 introduces the policy relevance of the present analysis and discusses the consequent methodological challenges in particular concerning the multiple nature of heterogeneity under investigation. Section 3 outlines the theoretical background of the farmers' behavioural responses to environmental measures on which the consequent research design and variable definition and selection are grounded. These latter are detailed in Section 4 after the presentation of the sample of Italian farms adopted in the study. Section 5 discusses the BCF estimation approach, highlighting its main advantages over other CML algorithms. Section 6 presents

the main results and assesses their robustness with respect to possible identification issues and biases. Section 7 concludes by highlighting some implications of the proposed methodology for the assessment and targeting of agri-environmental policies and sketches some directions for future research in the field.

## 2. Policy relevance and methodological challenges

Over the last few decades, the EU CAP has undergone several structural reforms and has increasingly emphasised the primary sector's environmental dimension (Commission of European Communities, 2000). Currently, the CAP includes objectives for protecting water, soil, climate and air quality, landscape, and biodiversity (European Commission, 2020). Following the 2014 CAP reform and the corresponding 2015–2020 CAP AEP design, these objectives are pursued via a diverse mix of policy instruments, three of which represent the subject of the present paper.

The oldest of these three means of intervention (introduced in 1992) consists of the AEMs. These are voluntary measures belonging to CAP's Pillar 2 which deliver compensatory payments to farmers to cover additional costs and foregone income from adopting more environmentally friendly practices. In our work with AEMs we refer to two measures which, after the 2014 CAP reform, are named "measure 10" (agri-environment-climate commitments) and "measure 11" (organic farming). Both measures provide monetary incentives for the voluntary adoption of eco-friendly farming techniques.<sup>ii</sup>

Following the 2003 "Agenda 2000" reform, a second environmental measure was introduced: CAP's Pillar 1 direct payments (DP) became subject to the so-called cross-compliance (CC) requirements that made these monetary subsidies contingent on several environmental and ecological standards. Although these requirements are intended to be mandatory, *strictu sensu*, complying with part of them is like satisfying an eligibility condition for first-pillar payments since non-compliance triggers administrative penalties up to the revocation of the DPs. Therefore, farmers may always give up applying for DP entirely, thus also ignoring part of the CC requirements.

The third policy instrument was introduced with the 2014 CAP reform through the so-called Greening Payment (GP). This measure represents the green' component of the new modified DP scheme, in which the financial support now hinges on three mandatory practices intended to benefit both the environment and the climate. Since it builds upon and reinforces CC, the GP is often regarded as a sort of additional (or super-) conditionality.<sup>iii</sup> As in the previous case, non-compliance results in a loss of support directly delivered to farmers. Therefore, under the 2014-2020 CAP design, eligibility for the full DP related to environmentally friendly practices now depends on satisfying both CC and the GP provisions.

It is worth noting that, in implementing such measures, there have been significant differences both across and within member states (MSs). For example, Italy has managed, implemented, and administered AEMs at the regional (NUTS-2) level through the so-called Rural Development Plans (RDPs). Similarly, although CC requirements have been enforced following the EU conditionality principles, the list of commitments applicable at the local level has also been left to the regional authorities. These include commitments to prevent soil erosion, organic matter decline, and soil compaction, perform a minimum level of ecosystem maintenance, and prevent habitat and landscape deterioration (National Rural Network, 2010). Finally, the GP is defined as a farm-specific, yearly, per-hectare payment calculated as a proportion of a farm's DP total value. Once again, however, the actual implementation of the GP may be differentiated at the regional level.

Therefore, MSs enforce and oversee these policy instruments acknowledging the existence of cross-country/cross-regional specificities, thereby allowing for some degree of flexibility in their implementation (Guerrero, 2021). Nevertheless, the content of all these intervention tools (i.e., their monetary implications and associated requirements) remains rigid in comparison to the very diverse conditions to which they apply. In fact, the same policy menu is offered to both very large farms and very small units, to both extensive livestock farming in mountain areas and orchards in plain urban areas, and so forth. This mismatch between highly heterogeneous farms and a relatively homogenous policy instrument is particularly delicate for Italy, whose primary sector mixes very different farming

traditions and peculiar geographical characteristics (Coderoni and Esposti, 2018). Such structural heterogeneity inevitably translates into behavioural heterogeneity in that the response of diverse farms to homogenous policies may substantially diverge in terms of both the size and nature of the response (i.e., the variables involved in the response). Moreover, even when farms exhibit analogous structural and behavioural characteristics, the uneven environmental effects that these policies may generate can result from very site-specific agronomic, ecological, and biophysical features such as field slopes, soil types, hydrology, and crop rotation (see for example Ó hUallacháin et al., 2016; Finn et al., 2009; OECD, 2022).

These multiple and complex sources of heterogeneity suggest that AEPs should be more flexible in targeting diverse farms. Unsurprisingly, the need for a more tailored design of the CAP environmental policies has frequently been advocated during the last two decades (Erjavec and Erjavec, 2015; Ehlers et al., 2021). In this respect, a policy rationalisation through better targeting of specific farm characteristics might help in achieving the declared environmental objectives, either through expenditure savings (for the same environmental performance) or improved environmental performance (for the same level of expenditure) (Esposti, 2022b). However, improving policy targeting and, ideally, tailoring, also requires a better understanding of whether and how the potential beneficiaries of such measures respond differently. Borrowing from the CI jargon, one would wish to identify and estimate HTEs (or individual TEs) as the natural empirical counterpart of this knowledge gap.

Policy evaluation studies addressing the impact of agri-environmental policies have gained considerable attention in recent years. Chabé-Ferret and Subervie (2013), Arata and Sckokai (2016), Mennig and Sauer (2020), and Bertoni et al. (2020), to name few recent examples, have applied difference-in-differences (DID) and/or matching techniques to assess the effects of different AEMs. Similarly, Bartolini et al. (2021) estimated the impact of AEMs in a multivariate treatment setting by adopting a generalised propensity score estimation. However, these studies typically have estimated average TEs (ATEs) without exploring TE heterogeneity, if not by focussing on specific farm groups

or considering quantile TEs (Esposti, 2017a; 2017b). The main risk of working with such aggregate measures is that of hiding systematically different unit or group-level effects. In other words, what holds true on average might not hold true for specific clusters and vice versa. Thus, may evidently lead to wrong policy conclusions.

In this respect, ML methods have recently proven a helpful toolbox for assessing AEPs. For example, Bertoni et al. (2021) used ML techniques to simulate the impact of GP in terms of land use change, although they did not touch on TE heterogeneity. Among the latest contributions, however, Stetter et al. (2022) represents the only study explicitly addressing the heterogeneous response of (South-eastern, German) farms to AEMs in terms of environmental performances. Interestingly, the paper acknowledges that the proper identification of such HTE can be problematic for at least two reasons: first, using the participation to AEMs as a binary treatment variable can only proxy for a wide range of sub-measures that farmers can choose from; second, measuring environmental performances is inherently hard because of the interconnected nature of many commonly adopted environmental indicators. Although HTEs can be particularly helpful for a better targeting of AEPs, thus improving their (cost-) effectiveness, these two caveats may complicate their empirical tractability.

On the one hand, when policy measures are delivered via sub-measures among which farmers can freely choose (i.e., a multi-valued treatment), the standard identification strategies for HTEs may fail due to the presence of alternative versions of the treatment (VanderWeele and Hernán, 2013; Lopez and Gutman, 2017). Moreover, the interpretation of the resulting estimand could be misguided because the local differences in TE could instead be driven by treatment heterogeneity (hereafter TH - Heiler and Knaus, 2022). On the other hand, had such disaggregation level been attainable, it would still be difficult to unambiguously link a specific scheme to a single environmental indicator. As previously mentioned, depending on both the farm's specificity and the treatment, elementary environmental outcomes are always interdependent and hard to examine in isolation (Chabé-Ferret and Subervie, 2013). In other words, for any treated unit, TEs can either differ across multiple indicators or, even worst, trigger spillovers such that changes in one environmental outcome may

impact others. Ignoring this output-dependent treatment effect heterogeneity (henceforth OTH) and focusing on elementary indicators may therefore lead to misleading interpretations of the HTE.

While our interest lies in estimating the HTE of both DPs and AEMs in general, this paper also acknowledges and attempts to empirically address the two issues discussed above.

### 3. Theoretical framework: modelling farmers' response to agri-environmental policies

We begin by discussing a simple theoretical framework conceptualising farmers' uptake of AEPs and providing a behavioural foundation for TE heterogeneity. Unlike the model presented in Stetter et al. (2022) where HTEs only result from farm-specific production technologies, we postulate a stylized behavioural mechanism explaining how farms respond to different policy option and, therefore, how HTEs may emerge. Moreover, our framework also formalizes how TH and OTH can interfere with the identification of the HTEs of interest.

Consider a panel of  $N$  production units (i.e., farms) observed over  $T$  time periods. Each farm can choose among  $K$  alternative AEPs. Next, assume that farmers are profit maximisers and, for simplicity, risk neutral. The latter greatly simplifies the following analytical treatment as it allows formulating farmers' behaviour in terms of actual profits ( $\pi_{it,k}$ ) rather than expected profits.<sup>iv</sup> In practice, we assume that none of the AEPs considered in this study implies a major change in the riskiness of farming activity.<sup>v</sup>

We postulate that each farm  $i \in \{1, \dots, N\}$  is associated with an aggregated general multi-input multi-output farm-specific technology represented by the feasible production set  $F_i \subset \mathbb{R}^M$ . Given  $F_i$ , the  $(M \times 1)$  vector of netputs  $\mathbf{y}_i = (y_{1i}, \dots, y_{Mi})'$  is feasible if  $\mathbf{y}_i \in F_i$ .<sup>vi</sup> This netput vector contains both farm-specific outputs (with positive signs) and farm-specific inputs' use (with negative signs), possibly including also non-market inputs and outputs. The adjective 'farm-specific' implies that  $F_i$  contains all possible sources of heterogeneity in the farmer's production decisions that depend on both external and internal factors (Esposti, 2022b).<sup>vii</sup> We can express the  $i^{th}$  farm's specific features with a  $Q$ -dimensional vector  $\mathbf{Z}_{it}$ .

To keep the notation consistent throughout the paper, we refer hereinafter to the set  $\{T_{it,1}, \dots, T_{it,K}\}$  as the treatment set and to  $T_{it,k}$  as treatment  $k$ . At period  $t \in \{0, \dots, T\}$ , any AEP chosen by farmer  $i$ ,  $T_{it,k}$ , is expected to induce specific production choices,  $\mathbf{y}_{it,k}$ , via either output production or input use. Therefore, treatments can be univocally mapped to production choices ( $T_{it,k} \leftrightarrow \mathbf{y}_{it,k}$ ). Notice that this argument also holds for multiple treatments. For example, suppose that the  $k$ -th treatment is delivered through  $V$  alternative versions ( $v = 1, \dots, V$ ) among which the farmers choosing the  $k$ -th treatment can choose (VanderWeele and Hernán, 2013). We can then indicate the treatment as  $T_{it,kv}$ . This does not affect the overarching structure of our theoretical model, as the new set of treatment option can be simply rewritten as  $\{T_{it,1}, \dots, T_{it,k1}, \dots, T_{it,kV}, \dots, T_{it,K}\}$ , and it is always possible to express ( $T_{it,kv} \leftrightarrow \mathbf{y}_{it,kv}$ ).

We can now express farmers' production choices as functions of the policy treatments themselves, given a farm-specific technology  $F_i$  as expressed in  $\mathbf{Z}_{it}$ , i.e.,  $\mathbf{y}_{it,k} = g(T_{it,k}, \mathbf{Z}_{it})$ , where  $g(\cdot)$  is a vector-valued function. In addition, if farms are profit maximisers and can choose  $T_{it,k}$ , the policy support operates like market price changes in orienting production decisions (Esposti 2017a; 2017b). Consequently, we can generically express farms' individual profit functions as  $\pi_{it,k} = \Pi[g(T_{it,k}, \mathbf{Z}_{it})]$ , where  $\Pi(\cdot)$  is a single-valued function.<sup>viii</sup>

This behavioural representation makes clear that farmers' choice is not driven by  $\mathbf{y}_{it,k}$ , which is the main target of the policy, but by the associated profit  $\pi_{it,k}$ . Following this logic, each observed pair  $(T_{it,k}, \mathbf{y}_{it,k})$  represents the profit-maximising combination of each treatment and the resulting set of production choices. Without assuming any specific functional form for the underlying technology or profit function, an augmented version of the weak axiom of profit maximisation can be formulated to identify the optimal netput vector  $\mathbf{y}_{it,k}$  (Afriat, 1972; Varian, 1984; Chavas and Cox, 1995; Esposti, 2000). This implies that  $\Pi[g(T_{it,k}, \mathbf{Z}_{it})] \geq \Pi[g(T_{it,h}, \mathbf{Z}_{it})], \forall k, h \in K, k \neq h$ . Namely, the profit of the  $i^{th}$  farmer choosing treatment  $k$  at time  $t$  ( $\pi_{it,k}$ ) exceeds the profit that she would have achieved

had she chosen any other alternative  $T_h$  ( $\pi_{it,h}$ ). For a given baseline treatment ( $T_{it,0}$ ), farm  $i$  will choose treatment  $k$  at time  $t$  if  $\Pi[\mathbf{y}_{it,k}(T_{it,k}, \mathbf{Z}_{it})] \geq \Pi[\mathbf{y}_{it,0}(T_{it,0}, \mathbf{Z}_{it})]$  or, alternatively,  $\Pi[\Delta g(T_{it,k}, T_{it,0}, \mathbf{Z}_{it})] \geq 0$ , where  $\Delta g = \Delta \mathbf{y}_{it,k} = \mathbf{y}_{it,k} - \mathbf{y}_{it,0}$ . Notice that, within this conceptual framework, the full treatment set might not be feasible for all farms. In fact,  $\mathbf{Z}_{it}$  might bind the choice of the netput vector  $\mathbf{y}_{it,k}$ , thereby limiting the choice of  $T_{it,k}$  to a subgroup of  $\{T_{it,1}, \dots, T_{it,K}\}$ . This may also apply when treatment is delivered through  $V$  alternative versions: given  $\mathbf{Z}_{it}$ , not all the sub-treatments,  $T_{it,k1}, \dots, T_{it,kV}$ , may be feasible for all farmers choosing the  $k^{\text{th}}$  treatment.

The main goal of this paper is to construct and identify an empirical counterpart of  $\Delta \mathbf{y}_{it,k}$  and determine its distribution across heterogeneous farms.<sup>ix</sup> Assuming that either  $\mathbf{y}_{it,k}$  or  $\mathbf{y}_{it,0}$  can be observed, this research question can be addressed using the CI analytical framework, where  $\Delta \mathbf{y}_{it,k}$  indicates the TE of interest, and  $\mathbf{y}_{it,0}$  represents the counterfactual state of  $\mathbf{y}_{it,k}$ , had the farm not chosen treatment  $k$  (Imbens and Rubin, 2015). However, in presence of multiple treatment versions ( $T_{it,k1}, \dots, T_{it,kv}, \dots, T_{it,kV}$ ),  $\Delta \mathbf{y}_{it,k}$  may differ from  $\Delta \mathbf{y}_{it,kv}$ , for some  $v \in V$ . In fact, not only may these two quantities differ but, more importantly, we may also observe  $(\Delta \mathbf{y}_{it,k} - \Delta \mathbf{y}_{jt,k}) \leq (\Delta \mathbf{y}_{it,kv} - \Delta \mathbf{y}_{jt,kv}) \neq (\Delta \mathbf{y}_{it,kv} - \Delta \mathbf{y}_{jt,kv})$  for any  $i, j \in N$  and any two  $v, v' \in V$ . Heiler and Knaus (2022) show that the above inequality results from  $(\Delta \mathbf{y}_{it,k} - \Delta \mathbf{y}_{jt,k})$  being a weighted average of all the treatment versions  $\Delta \mathbf{y}_{it,kv}$ , where the weights are proportional to the probability that farm  $i$  chooses  $T_{it,kv}$ . In other words, in presence of multiple treatment versions, we would erroneously mistake TE heterogeneity for what is, in fact, a diverse treatment choice mechanism, i.e., TH.

As introduced in Section 2, when it comes to evaluating the effect of a treatment, one could either focus on one or multiple elements of the netput vector  $\mathbf{y}_i = (y_{1i}, \dots, y_{mi}, \dots, y_{Mi})'$ . However, since most entries in  $\mathbf{y}_i$  can be highly interconnected (i.e.: some  $y$ s can be positively or negatively correlated with one or more other  $y$ s), evaluating TEs through marginal evaluations of these elements could make results hard to interpret. For example, consider any two positively (or negatively)

correlated items  $y_{mi}, y_{li} \in \mathbf{y}_i$ . Then, for any  $i, j \in N$  and treatment  $T_{it,k}$ , we will have that  $\Delta y_{mi,k}$  is also correlated with  $\Delta y_{li,k}$ . Therefore, comparing the marginal HTE for the two indicators, i.e.: comparing  $(\Delta y_{mi,k} - \Delta y_{mj,k})$  against  $(\Delta y_{li,k} - \Delta y_{lj,k})$ , can lead to misleading conclusions. We previously referred to this issue as OTH. In Section 4.3 we postulate that OTH can be addressed via dimension reduction, where we project a vector of correlated environmental indicators  $\mathbf{y}_i^e \subset \mathbf{y}_i$  onto a lower dimensional space through a synthetic environmental performance indicator. Nonetheless, it remains possible to empirically assess the potential interference of the OTH on HTE estimation by comparing the results obtained via the lower-dimensional index to those obtained on its individual components (See Section 6.3).<sup>x</sup>

If one can address TH and OTH, then under suitable restrictions on the joint distribution of the potential outcomes  $(\mathbf{y}_{it,k}, \mathbf{y}_{it,0})$  and given farm characteristics  $\mathbf{Z}_{it}$ , the identification of  $\Delta \mathbf{y}_{it,k}$  can be achieved via unconfoundedness (see Section 5) if  $\mathbf{Z}_{it}$  contains all the relevant variables that influence both the treatment choice,  $T_{it,k}$ , and the farmer's production choices (Angrist and Pischke, 2008; Wooldridge, 2010, Chapter 21; Imbens and Rubin, 2015, Chapter 3).

Following Brown et al. (2021) and Stetter et al. (2022), we distinguish between four sets of farm attributes:<sup>xi</sup> economic factors (i.e., factor endowment); socio-demographic characteristics (of the farm's holder and workforce); environmental (mostly geographical) factors; and idiosyncratic characteristics (of the farm's holder and workforce, such as ability, knowledge, motivations, beliefs, and values, as well as unobserved environmental features such as agronomic characteristics and fertility). To facilitate the illustration of our identification strategy, we assemble these characteristics into separate partitions of  $\mathbf{Z}_{it}$ , namely,  $\mathbf{Z}_{it} = (\mathbf{X}_{it}, u_i)$ , where  $\mathbf{X}_{it}$  consists of a  $(P \times 1)$  array. Furthermore, we define  $\mathbf{X}_{it} = (\mathbf{V}_{it}, \mathbf{S}_i)$ , where  $\mathbf{S}_i$  is a vector of observable time-invariant farm characteristics,  $\mathbf{V}_{it}$  is a vector of observable time-variant farm attributes.  $u_i$  represents unobservable time-invariant farm features. According to this categorisation, identifying HTEs requires two fundamental restrictions: first,  $\mathbf{V}_{it}$  must be pre-determined in that the treatment cannot affect  $\mathbf{y}_{it}$  via

$\mathbf{V}_{it}$ ; second,  $u_i$  must not be associated with both  $T_{it,k}$  and  $\mathbf{y}_{it}$ , under penalty of introducing selection-on-unobservable bias (Imbens and Rubin, 2015). Although the first condition can be satisfied using time-stable variables (i.e.,  $\mathbf{V}_{it} \approx \mathbf{V}_i$ ) or lagged values (see Section 4.4), the exogeneity of  $u_i$  is often assumed and tested via sensitivity analysis.

We maintain this assumption throughout the paper, thus only focusing on  $\mathbf{X}_{it}$  when discussing TE identification. As discussed in Sections 4.5 and 6.3, however, we also resort to suitable robustness checks to test the validity of our identification strategy under endogenous  $u_i$ .

## 4. Data and research design

### 4.1 Observational dataset

We use information from the Italian Farm Accountancy Data Network (FADN), which represents the only source of microeconomic agricultural data that is harmonized at EU level and collects physical, structural, economic, and financial data on farms in all EU Member States (European Council, 2009). The survey is representative of the farms that can be considered professional and market oriented, due to their economic size (that is equal or more than 8,000€ of standard output). In Italy these correspond to 95% of Utilised Agricultural Area, 97% of the value of Standard Production, 92% of Labour Units and 91% of Livestock Units. The representativeness of the dataset is ensured on three dimensions, namely: region, economic size, and farm typology. For these reasons the FADN is the most (and only) widely used farm-level dataset for, among others, CAP evaluations and specifically for the assessments of the AEP impacts (among others: Stetter et al., 2022; Bartolini et al., 2021; Arata and Sckokai, 2016).

Our research focuses on the 2014-2020 programming period of the CAP.<sup>xii</sup> However, unlike Stetter et al. (2022), we exclude the initial year (2014) for two reasons: first, payments of one of the policies under consideration (the GP) only started in 2015; second many of the farms observed in 2014 may still benefit from measures of the previous programming period. We thus focus on the 2015-2020 period, although we only have detailed and validated information until 2018. Therefore, our initial

sample consists of a representative collection of Italian commercial farms that produces an unbalanced panel consisting of 9,580, 10,135, 10,792, and 10,386 observations in 2015, 2016, 2017 and 2018, respectively. Because our analysis does not address regime-switching dynamics, we only consider farms for which the treatment status did not change over the period analysed, i.e.,  $T_{it,k} = T_{i,k}$  for all  $i \in \{1, \dots, N\}$ . For this reason, we first extract a balanced panel consisting of 5,836 units observed over the four-year period 2015–2018, then drop all entries satisfying  $T_{it,k} \neq T_{is,k}$  for any  $s, t \in \{2015, \dots, 2018\}$  and  $s \neq t$ .<sup>xiii</sup> The resulting dataset consists of 4,001 farms repeated over four years for a total of 16,004 observations. Compared to other related works (Bertoni et al., 2020; Stetter et al., 2022), our study provides wide coverage of the agricultural sector by focusing on the entire national area instead of a single region. Furthermore, since the treatments presented in Section 4.2 are likely to impact the agri-environment over several years, our outcome variable uses information from the last two years in the series to account for potential accumulation effects (see Section 4.3 for details).

#### 4.2 Definition of treatments

As mentioned in Section 2, the 2015–2020 CAP AEP design is based primarily on two main policy instruments that belong to either CAP’s Pillar 1, Pillar 2, or both. On the one hand, we observe Pillar 1 subsidies that are conditional on a set of compulsory requirements (i.e., CC and the GP) with which farmers must comply to preserve the DP. On the other hand, we have voluntary measures aimed at compensating farmers for income losses or increased costs resulting from the voluntary adoption of more sustainable farming practices (i.e., the AEM of Pillar 2). Consequently, farms are subscribed to – in fact, they voluntarily choose – one of three possible policy alternatives, which effectively reflect the interplay between the two pillars of the CAP: (i) farms failing to meet all the CC and GP requirements – that is, farms receiving neither Pillar 1 nor Pillar 2 payments; (ii) farmers receiving both Pillar 1 (DP and GP) and Pillar 2 (AEM) payments; and (iii) farms complying with the CC and GP requirements but not adopting any AEM.

Table 1 indicates how the farms in our sample are distributed across the three policy categories. The third cohort is the largest group, which includes approximately 71% of the observed farms (2,841 units). Using the terminology introduced in Section 3, we consider the corresponding policy option as the baseline treatment,  $T_{i,0}$ , associated with the netput vector  $\mathbf{y}_{it,0}$ . Next, all farms choosing not to benefit from Pillar 1 and Pillar 2 payments (i.e., the first cohort, corresponding to approximately 13% of the sample) take up the first treatment,  $T_{i,k=1}$ , which implies giving up both Pillar 1 and Pillar 2 resources. We assume that this decision follows the behavioural model stylised in Section 3, according to which, conditional on  $\mathbf{X}_{it}$ ,  $T_{i,1}$  produces higher profits than  $T_{i,0}$ . Similarly, farms applying for Pillar 2 AEM supports (i.e., the second cohort, corresponding to approximately 16% of the sample) choose treatment  $T_{i,k=2}$  through the same profit-maximising mechanism. In this respect, our work extends the analysis in Stetter et al. (2022) by distinguishing between the two different AEPs described above (i.e.: the AEMs and the Pillar 1 environmental requirements).

We postulate that treatments  $T_{i,1}$  and  $T_{i,2}$  belong to two non-overlapping choice sets, or in other words, we rule out a multiple treatment setup by positing treatment  $T_{i,1}$  as infeasible for farms choosing  $T_{i,2}$  and vice versa. Although this assumption is quite strong, it is necessary to identify the treatment effects of interest. However, given that  $T_{i,1}$  and  $T_{i,2}$  represent two ends of a rather wide spectrum of policy options, it is plausible that both treatments may appeal to (i.e., are feasible for) farms with very distinctive characteristics. Conversely, our setup implies that both  $T_{i,1}$  and  $T_{i,2}$  are feasible alternatives to the baseline treatment  $T_{i,0}$ . This presupposes that farms in the control group are characterised by features  $\mathbf{X}_{it}$  that overlap with the characteristic of the units in  $T_{i,1}$  or  $T_{i,2}$ . That is, we can always find comparable farms in either of the two groups within different strata of  $\mathbf{X}_{it}$ , i.e.,  $0 < Pr(T_{i,k} = 1 | \mathbf{X}_{it} = \mathbf{x}_{it}) < 1$ . This restriction is also commonly known as common support (or positivity), and as we discuss in Section 5 and Appendix E, it limits extrapolation issues, thus preventing unreliable TEs.

One caveat in our setup is that unlike  $T_{i,1}$ , farms choosing  $T_{i,2}$  may in fact opt for one among four treatment versions. As outlined in Section 2,  $T_{i,2}$  aggregates measure 10 and 11 which, in turn, can be decomposed in two sub-measures: Agri-environment-climate commitments (10.1); Conservation and sustainable use and development of genetic resources in agriculture (10.2); Payment to convert to organic farming practices and methods (11.1); Payment to maintain organic farming practices and methods (11.2). While measure 10.2 only concerns a small share of farms (roughly 3% of our sample) and can be thus excluded or safely merged into measure 10.1 (our current choice), sub-measures 11.1 and 11.2 are substantially equivalent in terms of farmers behaviour, the only difference being the amount of support granted. For this reason, we de-facto consider sub-measure 11.1 and 11.2 as a unique measure (i.e.: measure 11). As put forward in Section 2 and 3, disregarding such distinctions may greatly impact the interpretation of the HTEs via TH.

It is also worth mentioning that, in principle, the sub-measures discussed above could be further disaggregated into specific actions (using the RDP jargon). Unfortunately, the Italian FADN data do not provide enough information on AEM actions. In fact, to our knowledge, there are no high-quality representative datasets that can provide more detail on AEMs (see for example Stetter et al., 2022, who use the German version of our dataset). However, had this level of disaggregation been observable, it would imply a very large number of actions (i.e., treatment versions), as evidenced by the 21 RDPs implemented in Italy.<sup>xiv</sup> Clearly, expanding the treatment options well beyond the four sub-measures mentioned above would greatly affect the sample size of each subgroup and challenge the estimation of any HTE under the standard conditions discussed in Section 5 (Heiler and Knaus, 2022). Finally, focusing on more specific measures does not necessarily imply a more refined outcome variable (see Section 4.3 for further discussion).<sup>xv</sup>

Since organic farming (measure 11) is homogenous across the RDPs and involves a reasonable number of farms (271), we repeat our analysis by redefining treatment  $T_2$  as a two-versions treatment  $T_2 = (T_{2o}, T_{2n})$ , where  $o$  = organic and  $n$  = non-organic. Given our initial definition of the treatments,  $T_{2n}$  coincides with measure 10 which, unlike measure 11, is not entirely homogeneous across RDPs

and could thus be exposed to further TH. We therefore estimate the HTE of  $T_2$  under two different setups: (i) we first analyse the HTE of participating to AEMs as in Stetter et al. (2022); (ii) we then break down the treatment in (i) into  $T_{20}$  and  $T_{2n}$  and obtain the corresponding HTE. We finally compare the results from (i) and (ii) and discuss their implications for the interpretation of the HTE of interest (see Section 6.4).

**[Table 1 about here]**

### 4.3 The outcome variable

The theoretical framework presented in Section 3 expresses the farm response to the treatment as  $\Delta \mathbf{y}_{it,k}$ , that is, a vector whose non-zero elements represent all of the farmer's production choices associated with the treatment in terms of both input and output.<sup>xvi</sup> These elements may consist of a long list of the farmer's specific production decisions, ranging from crop and livestock management practices to water and nutrient use (Guerrero, 2021: 11; Burton and Schwarz, 2013). One way to reduce the dimensionality of  $\Delta \mathbf{y}_{it,k}$  consists of identifying and extracting the elementary indicators expressing the change in farming practices towards extensification or environmentally friendly practices. However, as discussed in Section 2 and 3, focussing on elementary indicators might cause ambiguity when interpreting TE heterogeneity because of the OTH problem. Given the potential correlation among the components of  $\mathbf{y}_{it,k}$ , one way to retain all the information in the netput vector while avoiding multiple marginal evaluations is to perform dimension reduction (Chipman and Gu, 2005) to obtain composite dimensional indices (Bartolini et al., 2021). Not only can this strategy provide an insulation against OTH, but it also resonates the need for a comprehensive evaluation of complex policy instruments such as the AEM discussed in Section 2 and 4.2. As also argued by Stetter et al. (2022: 727), despite the articulation of AEMs in specific sub-measures, the goal of the AEPs remains more general, aiming to improve the overall environmental performance of the agricultural sector. Although many studies have tried to evaluate the effectiveness of distinct AEPs with respect to specific policy targets (e.g., the impact on biodiversity), the integrated assessment of multifaceted

goals involving, for example, soil and water protection and the curbing of greenhouse gas (GHG) emissions have received relatively little attention until recently (Hudec et al., 2007; Zhen et al. 2022). However, the literature has long suggested that the intricate and ecosystemic nature of the agri-environment requires that any assessment should be based on a comprehensive integration of indicators across many environmental dimensions (Wascher, 2003; Purvis et al., 2009).

In this respect, Purvis et al. (2009) propose an interesting, harmonised approach to evaluating AEMs: the so-called Agri-environmental Footprint Index (AFI). The AFI expresses a multidimensional assessment as a univariate index that can be flexibly adapted to diverse contexts. We use the AFI framework as adapted by Westbury et al. (2011) with the FADN data. We refer to this methodology as FADN-AFI, as the resulting index uses elementary information included in the FADN dataset. We extend the FADN-AFI to evaluate whether and to what extent the implementation of the CC requirements, GPs and AEMs meet the CAP 2015-2020 environmental objectives.<sup>xvii</sup>

Table 2 presents the elementary components of our FADN-AFI<sup>xviii</sup>. The land use diversity indicator (the Shannon Index) is detailed in Appendix A. Appendix B discusses the definition of a farm-level GHG emissions indicator using farm-level information. This measure should provide a reliable proxy of the contribution of a farm's practices to climate change mitigation (Dabkiene et al., 2021). The FADN-AFI's elementary components are then standardised to obtain dimensionless z-scores that we eventually aggregate using the weights indicated in the last column of Table 2 (i.e., giving a positive or negative sign for positive or negative environmental externalities, respectively).<sup>xix</sup> The resulting FADN-AFI is monotonic in farms' environmental performance in that higher FADN-AFI scores correspond to 'better' environmental performance. Since the range of the FADN-AFI is not bounded, the index might be difficult to interpret per se. However, since HTEs are defined through pairwise differences, these can easily be understood comparatively. Finally, we average the FADN-AFI over years 2017–2018 to provide more stable values for the outcome variable.<sup>xx</sup>

**[Table 2 about here]**

#### 4.4 Confounding variables

As discussed in Section 3, the choice of covariates entering the  $\mathbf{X}_{it}$  vector becomes crucial for identifying the HTEs of interest. These should encompass farm heterogeneity as extensively as possible, thereby allowing fair comparisons between treated and untreated units. Selecting all the relevant confounders such that the assumptions outlined in Section 5 are satisfied may follow multiple routes. On the one hand, one may construct a very large collection of both internal farm characteristics and external socio-economic indicators that might explain the individual decision of adopting one of the treatments. In this case, we would let ML algorithm choose which of these features contributes the most to predict farmers' behaviour through a regularization mechanism. However, as recently outlined by Hünermund et al. (2023), this strategy may lead to severely biased TE if the covariate set includes potentially endogenous confounding variables. Ultimately, the authors advocate that, when the goal is conducting CI, researchers need to justify the controls that they wish including and, more importantly, make sure that these are exogenous (i.e.: pre-treatment).

For these reasons, we begin by defining the confounders in  $\mathbf{S}_i$  and  $\mathbf{V}_{it}$  through an extensive literature review covering several empirical studies addressing both farmers' participation in AEPs and the impact of AEPs on farms' economic and environmental performance. The results of this survey are displayed in Table 3, where the list of covariates resulting from this desk research is classified using the taxonomy elaborated by Brown et al. (2021) and discussed in Section 3. We invite the reader to refer to the individual studies for a throughout explanation of how these regressors are relevant for the two above mentioned research questions. The abundance of controls compiled in this long list might suggest some form of preliminary selection to avoid redundancy and achieve a more parsimonious set of variables. Nevertheless, unlike most parametric econometric tools, forest-based ML algorithms can easily accommodate multiple overlapping information sources and use them to either create intermediate features or discard redundant ones through regularisation. Therefore, our empirical analysis makes use of all the covariates in Table 3.<sup>xxi</sup>

To satisfy the identifying conditions anticipated in Section 3, however, the time-varying controls,  $\mathbf{V}_{it}$ , must be exogenous with respect to the treatment (i.e., pre-determined). In theory, this would preclude the use of certain direct measures of farm physical and economic size, such as utilised arable land, profit, revenue, costs, and total workforce. To circumvent this issue, some authors suggest using covariates measured before the introduction of the treatment (see, for example, Bertoni et al., 2020; Uehleke et al., 2022; Stetter et al., 2022 for studies assessing AEMs). However, this strategy is sometimes infeasible, as such measurements may not be available if the policies under investigation were introduced several years before the outcome is measured. When this happens, going back in time may imply a major loss of observations. This concern is particularly relevant for our application, as the rotating structure of the Italian FADN panel shows that 582 farms (approximately 15% of the sample) included in the 2015–2018 dataset are not present in the 2014 data. Therefore, our choice is to follow the strategy of Arata and Sckokai (2016) and Pufhal and Weiss (2009), which consists of using the first year since the introduction of the policy as the pre-treatment period<sup>xxii</sup> (2015, in this case). Notice that, since our outcome variable is calculated using the years 2017 and 2018,  $\mathbf{V}_{it}$  contains lagged (by 2 years) elementary components of the FADN-AFI. Moreover, since farms usually sign up for participating in certain AEMs over several years (Bertoni et al., 2020; Uehleke et al., 2022), we also include information on previous participation to such programs in  $\mathbf{V}_{it}$  (Chab  -Ferret and Subervie, 2013). Appendix Tables C1 and C2 report descriptive statistics for the outcome variable and all the control variables discussed above.

**[Table 3 about here]**

#### *4.5 Unobservable characteristics*

The theoretical derivation in Section 3 provides the behavioural foundation of the farmer’s treatment choice and response to the treatment. This behaviour depends on some observable characteristics but also on unobservable farm characteristics,  $u_i$ . The conditional independence between any of the treatments and the corresponding potential outcomes also hinges on the last component of the conditioning vector  $\mathbf{Z}_{it}$ , namely, the unobservable farm characteristics,  $u_i$ . If these latent features

influence both the choice between  $T_{i,1}$  and  $T_{i,2}$  and the corresponding potential outcomes, the identification of the HTE becomes challenging because of the violation of unconfoundedness. Even though  $\mathbf{X}_{it}$  can be extended to collect as many observable farm characteristics as possible, this strategy may be insufficient to insulate against selection-on-unobservable. Also, policy conclusions drawn from the HTE estimation could be problematic and even erroneous if the relevance of these unobservables and their possible association with the observable characteristics are not properly investigated and understood.

In these situations, ML methods (including BCFs) can help in the identification by automatically creating nonlinearities and complex interactions among the variables in  $\mathbf{X}_{it}$ , generating artificial strata that allow more precise comparisons between treated/untreated units and their counterfactuals. These ‘synthetic traits’ not only greatly expand the initial set of confounders but also correlate with the unobservable characteristics, thereby making the unconfoundedness assumption more credible. This argument is also put forward by Stetter et al. (2022: 738-739, 744), who provide a nice example of how this property of ML techniques may help to control for farmers’ attitudes towards environmental issues.<sup>xxiii</sup> Since this is not directly testable, we check the robustness of the above propositions through several sensitivity analysis tests. As illustrated in Appendix H, we probe the stability of our results in the presence of omitted variable bias from unobserved endogenous heterogeneity by introducing synthetically generated  $u_i$  into the covariate set. See Section 6.3 for more details and caveats of this approach.

## 5. Methodology

Research on the estimation of HTE has recently flourished, stimulated by an increasing interest in the development of ML methods able to provide theoretically sound inferences in such research settings (Athey and Imbens, 2019; Athey et al., 2019; Hahn et al., 2020; Knaus et al., 2020; 2021). Recent studies have proposed two ways ML can be used to estimate HTE. First, off-the-shelf ML algorithms can be tweaked to address some of the relevant identification issues of causal inference directly

(Athey and Imbens, 2016; Imai and Ratkovic, 2013; Wager and Athey, 2018; Hahn et al., 2020).<sup>xxiv</sup>

Second, direct modifications of the loss functions and data-splitting techniques can also help address one challenging problem of traditional ML techniques in causal settings: regularisation-induced confounding (hereinafter RIC - Nie and Wager, 2021; Hahn et al., 2018; Hahn et al., 2020; Chernozhukov et al., 2018, and references therein). We broadly refer to all these methods as CML.

Among the diverse approaches proposed in the literature, BART-based algorithms (Chipman et al., 2010; Hill, 2011; Hill et al., 2020) stand out as promising additions to the CML toolbox. These methods not only exhibit encouraging performance in terms of unbiasedness and coverage rates (Dorie et al., 2019; Carvalho et al., 2019; Hanh et al., 2020; Lee et al., 2020), but also take advantage of a fully probabilistic (i.e., Bayesian) inferential approach which enables the introduction of uncertainty measures when performing comparisons between groups of individuals (an aspect that currently limits the extent of other comparable ML methods - Stetter et al., 2022) and facilitates investigating the extent of overlap between treated and untreated groups (Hill and Su, 2013; Li et al. 2022) (see Appendix E for details). The latter is particularly important when it comes to treatment  $T_1$ , as the farms associated with this group are likely to exhibit very specific characteristics (see Appendix C and Esposti, 2017a; 2017b). Both traits hinge on the full posterior distributions of, on the one hand, the estimated HTE and, on the other hand, the fitted individual-level conditional expectations.

As with many other tree-based methods, BART can flexibly fit complex response surfaces by creating regularised ensembles of shallow Bayesian regression trees (Chipman et al., 1998), making it possible to perform predictive inference using the resulting posterior distributions (Chipman et al., 2010). This flexibility is achieved via recursive partitioning of the covariate space at the tree level, a procedure that is adept at defining nonlinearities and interactions between the observed covariates without the need to pre-specify them (Hill, 2011). However, since the original BART was not purposely designed for CI, a naïve application of such methods for the estimation of HTE might potentially introduce RIC. For this reason, Hahn et al. (2020) have recently proposed an extension of the original algorithm, which they refer to as Bayesian causal forests<sup>xxv</sup> (BCF). In addition to exploiting the estimated

propensity score (PS) to deal with potential distortions attributable to RIC (see Appendix D), the BCF algorithm also provides for a more flexible structure which separates the prognostic component from the heterogeneous treatment effect, thereby enabling direct control over the latter to avoid overfitting.

### 5.1 Estimating treatment effects via BCF

The estimation of HTEs using the BCF algorithm requires the usual assumptions of unconfoundedness and SUTVA, which can be expressed as follows:

$$Y_i(0), Y_i(1) \perp T_{i,k} | \mathbf{X}_i \quad (1)$$

where  $Y_i$  represents the FADN-AFI defined in Section 4.3,  $\mathbf{X}_i$  indicates the vector of confounders defined in Section 4.4, while  $Y_i(1)$  and  $Y_i(0)$  indicate potential outcomes for individuals in a treatment group ( $T_{i,k} = 1$ ) or control group ( $T_{i,k} = 0$ ), respectively (Imbens and Rubin, 2015, chapter 1). Notice that SUTVA implies no hidden variations of the treatment. Therefore, as discussed in Sections 2 and 3, binarized multiple-versions treatments can lead to violations of this assumption unless one imposes stringent restrictions on the treatment assignment mechanism. For example, in case any individual  $i$  with characteristics  $\mathbf{X}_i$  can only choose one of the hidden treatments, SUTVA is still a credible assumption (VanderWeele and Hernan, 2013; Lopez and Gutman, 2017). As previously discussed, we make this assumption for the treatments defined in Section 4.2, except for the distinction between organic and non-organic farming. We therefore set  $k$  to  $k \in \{1,2\}$  such that  $T_{i,k} = 1$  indicates either  $T_{i,1} = 1$  or  $T_{i,2} = 1$ , while  $T_{i,k} = 0$  always refers to farms in the control group. We discuss the implication for disaggregating  $T_{i,2}$  into  $T_{i,2o}$  and  $T_{i,2n}$  in section 6.4. For notational convenience, we henceforth drop the subscript  $k$ . Of these elements, we only observe the potential outcome that corresponds to the realised  $T_i$ , namely,  $Y_i = T_i Y_i(1) + (1 - T_i) Y_i(0)$ . Equation (1) postulates independence between the potential outcomes and the treatment, conditional on the set of exogenous variables,  $\mathbf{X}_i$ .

Combining unconfoundedness, SUTVA, and overlap (as discussed in Section 4.2) allows the estimation of causal effects via strong ignorability, i.e.,  $\mathbb{E}[Y_i(t)|\mathbf{X}_i = \mathbf{x}_i] = \mathbb{E}[Y_i|T_i = t_i, \mathbf{X}_i = \mathbf{x}_i]$ , with  $t_i \in \{0,1\}$ . The latter implies that the estimand of interest is simply the difference between two conditional expectation functions:

$$\begin{aligned}\tau(\mathbf{x}_i) &= \mathbb{E}[Y_i|T_i = 1, \mathbf{X}_i = \mathbf{x}_i] - \mathbb{E}[Y_i|T_i = 0, \mathbf{X}_i = \mathbf{x}_i] \\ &= \mu_1(\mathbf{x}_i) - \mu_0(\mathbf{x}_i)\end{aligned}\tag{2}$$

where  $\tau(\mathbf{x}_i)$  is typically referred to as a conditional average treatment effect (CATE). Since one can use  $\mu_T(\mathbf{x}_i)$  to impute conditional treatment effects at the individual level, (2) is sometimes referred to as individualised average treatment effect (IATE) (Lechner, 2018; Knaus et al., 2020; 2021). This estimand represents most disaggregated form of HTE.

Often, however, researchers may be interested in subgroups or intermediate aggregation levels of the exogenous covariates, leading to the definition of group average treatment effects (GATEs):

$$\tau(\mathbf{g}_i) = \int d\mathbf{x}_i \tau(\mathbf{x}_i) \phi_{\mathbf{X}_i|\mathbf{G}_i=\mathbf{g}_i}(\mathbf{x}_i)\tag{3}$$

where  $\phi(\cdot)$  represents a generic probability density of mass function,  $\mathbf{G}_i$  denotes the collection of possible groups, and  $\mathbf{g}_i$  denotes one such group. GATEs have recently gained considerable attention in the applied literature as treatment effect heterogeneity is often better understood for subsets of the population (Lechner, 2018; Lee et al., 2020). Average treatment effects (ATEs) can also be obtained by averaging the IATEs over the full distribution of  $\mathbf{X}_i$ :

$$\tau = \int d\mathbf{x}_i \tau(\mathbf{x}_i) \phi_{\mathbf{X}_i}(\mathbf{x}_i)\tag{4}$$

To estimate the IATEs (and then the GATEs and ATEs), we assume that the data-generating process for  $Y_i$  follows a stochastic process defined as follows:

$$Y_i = f(\mathbf{X}_i, T_i) + \varepsilon_i\tag{5}$$

where  $f$  indicates an arbitrarily complex function<sup>xxvii</sup> and  $\varepsilon_i$  represents an additive idiosyncratic error term  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , independently distributed.

In this context,  $\mathbb{E}[Y_i|T_i = t_i, \mathbf{X}_i = \mathbf{x}_i] = f(\mathbf{x}_i, t_i)$  therefore, at least in principle,  $\tau(\mathbf{x}_i)$  can be estimated by the simple difference  $f(\mathbf{x}_i, t_i = 1) - f(\mathbf{x}_i, t_i = 0) = \mu_1(\mathbf{x}_i) - \mu_0(\mathbf{x}_i)$ , as illustrated above. However, as discussed by Künzel et al. (2019) and Nie and Wager (2021), training two separate conditional mean functions and taking their difference may produce highly unstable estimates. For this reason, Hahn et al. (2020) proposed a slightly different approach, wherein the expected value of the outcome of interest has two components: a prognostic function,  $m(\mathbf{x}_i)$ , plus an additive heterogeneous treatment effect,  $\tau(\mathbf{x}_i)$ :

$$\mathbb{E}[Y_i|T_i = t_i, \mathbf{X}_i = \mathbf{x}_i] = f(\mathbf{x}_i, t_i) = m(\mathbf{x}_i) + \tau(\mathbf{x}_i)t_i \quad (6)$$

where both  $m(\cdot)$  and  $\tau(\cdot)$  represent stochastic functions with BART priors, namely:  $m \sim \text{BART}(\boldsymbol{\theta}|\widehat{\text{PS}}(\mathbf{x}_i), \mathbf{x}_i)$  and  $\tau \sim \text{BART}(\boldsymbol{\vartheta}|\mathbf{x}_i)$ , and  $\widehat{\text{PS}}(\mathbf{x}_i)$  indicates the estimated PS. The two vectors  $\boldsymbol{\theta}$  and  $\boldsymbol{\vartheta}$  collect the hyperparameters regulating the number of trees in the BART ensembles, their depth, and the splitting rule associated with each single tree (see Appendix F for details). As previously mentioned, the specification in Equation (6) allows regularising  $\tau(\mathbf{x}_i)$  directly and independently, thereby reducing the noisiness of the IATEs with respect to the same estimates obtained from simple differences in conditional mean functions. Furthermore, the additive nature of Equation (6) ensures that the prior on  $f(\mathbf{x}_i, t_i)$  is also a BART (Chipman et al., 2010; Hill et al., 2020). Finally, notice that the model presented in Equation (6) also appears in Nie and Wager (2021), who propose a frequentist approach to estimating  $\tau(\mathbf{x}_i)$ . In contrast to the setup discussed above, however, the authors propose a residuals-on-residuals re-parametrisation of Equation (6) which is then used to obtain (regularised) consistent estimates of  $\tau(\mathbf{x}_i)$  via a two-stage optimisation procedure. The full Bayesian model requires the definition of a likelihood function for the outcome variable (Gelman et al., 2013; McElreath, 2020). Consistent with Equation (5) and Chipman et al. (2010), Hill (2011), and Hahn et al. (2020), we employ a normal model for  $Y_i$ , along with a semi-conjugate inverse-Chi square prior for its variance:

$$Y_i \sim \text{Normal}(m(\mathbf{x}_i) + \tau(\mathbf{x}_i)t_i, \sigma^2) \quad (7)$$

$$m \sim \text{BART}(\boldsymbol{\theta} | \widehat{\text{PS}}(\mathbf{x}_i), \mathbf{x}_i)$$

$$\tau \sim \text{BART}(\boldsymbol{\vartheta} | \mathbf{x}_i)$$

$$\sigma^2 \sim \text{Inv}\chi^2(\omega)$$

where  $\omega$  is set following Chipman et al. (2010) (see Appendix F for further details). Samples from the posterior distribution of  $\tau(\mathbf{x}_i)$  are obtained via Markov chain Monte Carlo (MCMC) sampling, as implemented in the R package `bconf`. We indicate posterior draws from  $\phi(\tau(\mathbf{x}_i) | \mathbf{x}_i, t_i, y_1, \dots, y_N)$  as  $\{\tau^s(\mathbf{x}_i)\}_{s=1}^S$ , where  $S$  indicates the number of MCMC simulations.

### 5.2 Subgroup search via shallow regression trees.

The approximated posterior  $\{\tau^s(\mathbf{x}_i)\}_{s=1}^S$  is a multivariate probability distribution over a complex  $P$ -dimensional function, and as such, it might be difficult to interpret directly. One way to compress such information consists of obtaining marginal distributions of the IATEs for one covariate of interest and plotting them against the full range of that variable. A similar approach was adopted by Stetter et al. (2022), who used Shapely values (Shapley, 1953) to identify the marginal contributions of several treatment effect drivers and employed these indicators to construct partial dependence plots. Another sensible approach to investigating IATE heterogeneity consists of comparing farm subgroups obtained by projecting the full posterior distribution onto a lower-dimensional covariate space. In this respect, we follow the work of Hahn et al. (2020), Woody et al. (2021) and Yeager et al. (2019) (and, partially, Lee et al., 2020), who suggest eliciting the relevant subgroups by partitioning the IATE maximum-a-posteriori (MAP) estimates,  $\check{\tau}_i = S^{-1} \sum_{s=1}^S \tau_q^s(\mathbf{x}_i)$ , using shallow regression trees (CART – Breiman et al., 1984). Specifically, the authors propose to split  $\check{\tau}_i$  along  $\mathbf{w}_i$ , where  $\mathbf{w}_i \subseteq \mathbf{x}_i$  indicates a vector of policy-relevant variables and setting  $\mathbf{w}_i \subset \mathbf{x}_i$  implies using domain knowledge to enforce an initial regularisation of the resulting tree. In this paper, we restrict our attention to a subset of simple and understandable characteristics that policymakers might find helpful to improve the targeting of AEMs (see Section 6.2). Once farm subgroups have been identified, GATEs can be obtained as weighted averages of the IATEs that fall into each cluster. This

approach to calculating GATEs is also consistent with Lechner (2018) in that group-level effects are obtained as convex combinations of the IATEs. In our application, however, weighting is automatically performed when fitting a tree to  $\check{\tau}_i$ .

Finally, for some potential effect moderator  $x_p \in \mathbf{x}$ , the comparison between pointwise estimates (or intervals) computed at different levels of  $x_p$  ignores any potential correlation between IATEs along other variables  $x_l$ , for all  $p, l \in \{1, \dots, P\}$ . In other words, the marginal distribution of  $\tau(x_{p,i})$  disregards the information encoded in the correlation between  $\tau(x_{l,i})$  and  $\tau(x_{l,-i})$  when  $x_{l,i}$  and  $x_{l,-i}$  are close. This might lead to misleading comparisons along  $x_p$  and, consequently, unreliable policy implications. Therefore, once the relevant subgroups have been identified, one can obtain the full posterior distribution of each pairwise difference as:  $\phi_{g_1, g_2} = \phi(\tau_{i|i \in g_1} - \tau_{i|i \in g_2})$ , where  $g_1$  and  $g_2$  indicate any two subsets of  $\check{\tau}_i$ .

## 6. Results

### 6.1 IATEs

Figure 1.1 (both A and B) displays the MAP (i.e., the average over the  $S$  samples from the posterior distribution of  $\tau(\mathbf{x}_i)$ ) estimates and corresponding 95% confidence intervals (CrI) of the IATEs over the two treatment comparisons. These are ordered across the respective samples from the lowest to the highest individual value. We start our discussion by presenting the results for  $T_2$ , the treatment that is more frequently addressed by the literature. First, it is worth noting that, overall, the modal direction of the responses to the treatment ( $T_2$ ) is fully consistent with theoretical expectations: adding the AEM to the environmental standards implied by the CC and the GP (Figure 1.1-A) induces an improvement in the FADN-AFI, that is, in the farm-level environmental performance. The opposite response is observed when the environmental standards implied by the CC and GP are dropped (i.e., treatment  $T_1$ ) (Figure 1.1-B). However, whereas in the first case, most estimated IATEs exhibit CrI not including zero (black dots), the converse applies to the second comparison group, for which a

large proportion of farms have inconclusive individual-level TEs (light grey dots). Figure 1.1-B also indicates that some farms might even exhibit opposite responses, although the corresponding ITAEs appear quite noisy. This evidence is presented in greater detail in Table 4, which provides descriptive summaries of our main results.

Figure 1.2 (both A and B) show the IATE's MAP frequency distribution for the two cases. These plots highlight the variability of the responses, with few cases showing a treatment effect direction that conflicts with the expected direction (despite exhibiting CrI including both positive and negative values). Apart from these rare extreme cases, however, our MAP estimates range between roughly 0.1 and 1.0 for treatment  $T_2$  and between approximately -3 and 1.5 for treatment  $T_1$ . The nature and determinants of these different patterns can be further investigated by estimating GATEs, as addressed in the next section.

The irregularity of farms' responses to the treatments is a clear sign of heterogeneity, one that would be lost by the mere inspection of ATEs (see Figure 1.3, both A and B). Whereas these latter aggregated estimands provide clear indications of policy effectiveness (as both show an effect in the expected direction), the inspection of the IATEs tells a different and more subtle story. This is especially true for the treatment  $T_1$ , whereas the responses seem more homogeneous when studying the TE of implementing CC and GP requirements together with AEMs (treatment  $T_2$ ).

**Figure 1.** Estimated IATEs for treatment groups  $T_2$  and  $T_1$ . 1: point estimates of IATEs (median line) ordered from the smallest to the largest. The upper and lower dots represent the posterior 95% CrI endpoints associated with each individual MAP point estimate. 2: distribution of the MAP point estimate displayed in panel 1. 3: Monte Carlo approximation of the full posterior distribution of ATE defined as in Equation (4).

**[Figure 1 about here]**

**[Table 4 about here]**

Finally, for each individualised TE, we calculate the posterior probability that the corresponding IATE is either greater than zero or lower than zero for the  $T_2$  and  $T_1$ , respectively. Our results show that, when comparing farms implementing CC and GP requirements plus AEMs with the control

group, most of the IATEs' posterior distributions lie above zero. For example, the proportions of IATEs with at least 60%, 75%, and 90% positive posterior are 100%, 88.5% and 5%, respectively. Conversely, when comparing the control group to farms with no adherence to either CC or GP, the posterior distributions of their IATEs are largely negative. In this case, the proportions of IATEs with at least 60%, 75%, and 90% negative posterior are 83%, 15%, and 0%, respectively.

Notice that all the results discussed thus far are based on observations satisfying the common support as defined by Rule I in Appendix E. Under such a restriction to the range of  $\mathbf{X}_i$ , however, our dataset does not suffer drops. The sensitivity of these figures to different exclusion rules is discussed in Section 6.3 (Robustness check), in which the selection method we employed based on the estimated PS is also discussed.

## 6.2 GATEs

We partition the posterior distribution of  $\tau(\mathbf{x}_i)$  using a set of policy-relevant measures  $\mathbf{w}_i$  covering the most relevant dimensions of heterogeneity, as evidenced by the measures of feature importance produced by the BCF. Our characterisation of  $\mathbf{w}_i$  involves (i) examining the variable importance metrics generated as a by-product of the fitting model (7)<sup>xxviii</sup> and (ii) choosing the ten most predictive dimensions that policymakers might target to improve the effectiveness of AEMs. We next fit a CART algorithm to  $\check{\tau}_i$  using the attributes selected using the procedure illustrated above: latitude, longitude, altitude (geographical location); total arable land, share of rented land, revenue (physical or economic size); farm specialisation (relative importance of the first and second crop, farms specialised in livestock, crop and livestock farms, farms specialised in annual crops, and farm specialised in perennial crops). The results for the two treatments are shown in Figures 4 and 5, wherein, for the sake of interpretability, we do not allow the trees to split more than three times.

When we consider the adoption an AEM in addition to CC and GP requirements (treatment  $T_2$ ) (Figure 2.1-B), we find that TE heterogeneity is mostly associated with five variables: latitude, physical farm size, altitude, crop specialisation (share of the second crop in the crop mix), and

livestock intensity. These covariates trace out eight subgroups with different levels of TE. For example, subgroup  $g_8$  exhibits the lowest TE and consists of farms in southern Italy with less than 85 hectares of arable land. On the opposite end of the spectrum, we find subgroup  $g_{15}$ , which comprises crop-specialised farms in northern Italy with low livestock intensity. One can then obtain the full posterior distribution of  $g_{15} - g_8$  with 95% CrI between -0.27 and 0.49 (Figure 2.2-B), which indicates that the difference between the two subgroups is in fact small, if not zero. Interestingly, if we repeat this exercise across all the leaves defined by the tree in Figure 2.1-B, no group differences emerge (see Appendix G). These results are consistent with our discussion in Section 6.2, i.e., our preliminary findings suggested limited TE heterogeneity for treatment  $T_2$ .

In the case of treatment  $T_1$  (Figure 2.1-A), we see that the shallow tree picks up four moderating variables: specialisation in perennial crops, latitude, altitude, and livestock intensity. In this case, the subgroup with the strongest TE is  $g_8$ , which consists of farms specialised in perennial crops in Italy's southmost regions. Subgroup  $g_{15}$  includes observations from farms in the Po Valley that are not specialised in perennial crops. The difference in TE between these subgroups lies approximately between -2.2 and -0.41 (95% CrI – Figure 2.2-A), indicating the presence of TE heterogeneity. Repeating this exercise across all the terminal nodes, we find that, unlike treatment  $T_2$ , when the treatment consists of dropping both CC and GP requirements, many groups exhibit diversified responses. These further details are provided in Appendix G, where we also provide a deeper tree to gain further insights into these HTEs and a graphic representation of the geographical distribution of the IATEs.

*Figure 2. Shallow regression tree fitted to the MAP IATEs.*

**[Figure 2 about here]**

It is finally worth stressing that, although our main goal is to explore which observable farm characteristics exhibit a greater heterogeneity of response, some of these features might not be easily addressed by AEPs due to cost constraints or infeasibility, or because they could potentially lead to

discriminatory outcomes. Therefore, from a policy perspective, it would be more useful to evaluate the level of heterogeneity associated with covariates that can be targeted more easily and effectively through policy measures. In particular, most of the geographical features considered in our study, along with variables indicating long-term farm production specialization, appear particularly suitable for this purpose. In this respect, our results confirm that most of these geographical features significantly contribute to the observed heterogeneity of response. Similarly, the presence of perennial crops, crop specialization and livestock density, all of which relate to distinct and consistent farming practices, pinpoint to patterns of strong heterogeneity. This suggests that AEPs could significantly enhance their effectiveness by specifically targeting these features. For a more detailed discussion on this matter, please refer to Appendix G.

### *6.3 Robustness checks*

In this section, we check the consistency of our results to the assumptions formulated in Sections 4.4, 4.5, and 5.1. Our first robustness check concerns the common support condition. As anticipated in Section 5 and further detailed in Appendix E, we use both the posterior distribution of the BART algorithm and a PS-based algorithm to investigate common support. Our tests show that the results presented in Sections 6.1 and 6.2 are robust to these different methods to achieve overlap (see Appendix Tables H1 and H2).

We then perform a battery of tests that largely encompass those discussed by Stetter et al. (2022) in that we re-estimate our BCF multiple times, each time manipulating different model features. We begin by probing unconfoundedness through a recursive procedure in which we fit model (7) after dropping: (i) the most important feature in terms of relative frequency within the forest; (ii) the three most important features, and (iii) the five most important features. As detailed in Appendix Figures H1–H3), this exercise yields the first indication that the BCF in Equation (7) is fundamentally resilient against unobserved heterogeneity as long as this is associated with the set of observed confounders. Put differently, the complex interactions and nonlinearities generated by the tree ensemble seem to

work as additional synthetic controls associated with the left-out covariates, thus compensating for their absence in the model. However, this line of reasoning hinges on the (strong) assumption that the most predictive features are also associated with both  $Y$  and  $T_k$ . In case this assumption fails, the procedure discussed above cannot be interpreted as a robustness check for unconfoundedness. For this reason, we build upon these preliminary results and devise an additional test targeting endogenous unobserved heterogeneity directly. Our strategy consists of generating a random variable correlated with both  $Y$  and  $T_k$ , forming the vector  $\mathbf{Z}_{it}$  as described in Section 3, and re-running the model. As shown in Appendix Figure H4, our results do not change substantially, even under a strong imposed association between the unobserved variable and  $(Y, T_k)$ . This stability could result from the properties of the BART ensemble in that, when the forest is dense, the marginal contribution of each covariate becomes increasingly small (Chipman et al. 2010). Alternatively, it could be that the correlation between the nonlinear interactions generated by the BCF and the new confounder is strong enough to prevent distortions in the IATEs. In either case, it is worth warning that TE estimates might deteriorate quickly when unobserved heterogeneity is more abundant and complex. This test is in fact only restricted to a single unobserved factor, which we model as linearly associated with  $Y$  and  $T_k$  (i.e.: through correlations, which do not necessarily imply a direct effect of the synthetic  $u_i$  on either the outcome or the treatment). We thus expect that in presence of multiple endogenous latent confounders, possibly related to the treatment and the outcome (or other elements of  $\mathbf{X}_i$ ) in a non-linear fashion, our estimates might turn out sensibly different. Although the literature offers other methods to perform sensitivity analysis with respect to omitted confounders (Dorie et al., 2016; VanderWeele and Ding, 2018), we believe that they either do not overcome the limitations discussed above, or they remain difficult to implement in HTE estimation. Therefore, despite the promising results presented so far, we stress that these only hold if several important restrictions are met.

The following robustness check consists of creating both a placebo treatment and a placebo outcome, replacing their observed counterparts in Equation (7), and fitting the model two more times. If the model is correctly specified, the IATEs resulting from these ‘fake’ variables should be uncorrelated

with  $\tau(\mathbf{x}_i)$ . As Appendix Figures H5 and H6 show, the new results obtained through placebo treatments and outcomes not only have no correlation with our estimated IATES but also produce zero ATE with minimal TE heterogeneity.

Finally, we assess the robustness of the estimated IATES with respect to the OTH problem discussed in previous sections. We proceed by replacing the FADN-AFI with its elementary components and re-estimating model (7) as discussed throughout the paper. Appendix Figure H7 suggests that focussing on marginal indicators produces TEs whose individual directions are essentially in line with those presented Section 6.1. For example, implementing AEMs seems to yield lower GHG, higher crop diversity, lower fertilizer expenditure and more woodland areas. Nonetheless, a noteworthy difference emerges in terms of TE heterogeneity. Whereas adopting the FADN-AFI points to a limited diversity across farms, using marginal measurements would suggest that treatment  $T_2$  is environmentally beneficial only when the TE is large. For this reason, our results invite to caution when it comes to choosing the dependent variable of model (7). Although addressing individual indicators may appear more attractive and interpretable, it is worth stressing that missing out on the potential correlation or interdependence among them can affect the TE estimates in a nontrivial way.

#### *6.4 The role of heterogeneous treatments*

As discussed in Section 3 and 4, one potential limitation of our results (as well as other works investigating HTE of aggregated treatments) is that part of the estimated treatment effect heterogeneity in  $T_2$  might be a statistical artefact. This would result from the fact that  $T_2$  is a multiple-versions treatment as it aggregates two distinct measures which admit, in turn, several sub-measures (see Section 4.2). As introduced in Section 5, the presence of TH may affect our results by violating SUTVA (Heiler and Knaus, 2022). Since, in this case, the resulting interpretation of  $\tau(\mathbf{x})$  would be misleading, we re-estimate model (7) replacing  $T_2$  with the two respective measures (measure 11 and measure 10) and approach the problem from a multiple-versions treatment perspective as discussed in Lopez and Gutman (2017).

In order to assess the possible bias in HTE estimation due to TH, we compare the posterior distribution of the IATEs presented in Section 6.1 with the posterior density of the IATEs estimated using either  $T_{2o}$ ,  $\tau_{2o}(\mathbf{x}_i)$ , or  $T_{2n}$ ,  $\tau_{2n}(\mathbf{x}_i)$ . Figure 3 shows the 95% CrI for the differences  $\tau_{2o}(\mathbf{x}_i) - \tau_2(\mathbf{x}_i)$  and  $\tau_{2n}(\mathbf{x}_i) - \tau_2(\mathbf{x}_i)$ , respectively, where  $\tau_2(\mathbf{x}_i)$  indicates the IATE for individual  $i$  under treatment  $T_2$ . As we can see from these plots, the difference between our initial estimates and those obtained by substituting  $T_2$  with  $T_{2o}$  are minimal. Indeed, although  $\tau_{2o}(\mathbf{x}_i)$  is on average (black line in Figure 3.1) slightly smaller than  $\tau_2(\mathbf{x}_i)$  for all  $i \in N_o$ , where  $N_o$  indicates the number of units choosing  $T_{2o}$ , all the CrI include both positive and negative values. At the same time, when focussing on  $T_{2n}$ , we see that  $\tau_{2n}(\mathbf{x}_i) - \tau_2(\mathbf{x}_i)$  are on average higher than zero for all  $i \in N_n$ , where  $N_n$  indicates the units choosing  $T_{2n}$ . However, the CrI once again include zero for all such comparisons, although they are all moderately skewed towards positive values. Moreover, as mentioned in Section 4.2,  $T_{2n}$  could still entail some degree of TH, which recommends caution when interpreting the corresponding estimates. Overall, examining the two measures separately highlights that the posterior distribution of the IATEs does not seem to change markedly when the aggregated ( $T_2$ ) or the disaggregated ( $T_{2o}$ ,  $T_{2n}$ ) treatment is considered. This would suggest a limited impact of TH on our interpretation of the HTEs discussed above. Nonetheless, further research effort remains desirable to better clarify the possible role of multiple versions in the correct identification and estimation of the HTE.

**Figure 3.** IATEs differentials for the two versions of  $T_2$ . 1: 95% CrI for the distribution of  $\tau_{2o}(\mathbf{x}_i) - \tau_2(\mathbf{x}_i)$  (grey lines) and corresponding posterior means (black line). 2: 95% CrI for the distribution of  $\tau_{2n}(\mathbf{x}_i) - \tau_2(\mathbf{x}_i)$  (grey lines) and corresponding posterior means (black line).

[Figure 3 about here]

## 7. Concluding remarks

Giving the CAP a more explicit environmental orientation and justification has been at the core of all its recent reforms. This necessarily means shifting the support from undifferentiated and unconditional payments to more tailored and target measures. The efficiency and effectiveness of

agri-environmental policies in this respect critically depend on how farmers respond to these measures. This response, in turn, largely depends on the individual characteristics of supported farms. This makes the response itself highly heterogeneous and, consequently, suggests that there is still room for substantial improvement through better policy targeting.

In this paper, we present a causal machine learning approach to assess the heterogeneous response of farmers to different agri-environmental policies implemented through the 2015-2020 CAP reform. Building on the existing literature, our work's main contribution is twofold. Firstly, we explicitly conceptualize and investigate the different sources of heterogeneity that we expect influence farms' environmental performances under such policies. Secondly, we take advantage of the most recent developments in Bayesian non-parametrics and conduct the analysis using a relatively unexplored algorithm called Bayesian causal forest. This method allows using the posterior distribution of the individualised TE (i.e., the IATEs) to draw inferences about arbitrary transformations of these highly disaggregated estimands. We leverage this property throughout the paper, particularly when discussing group-level treatment effects and testing the robustness of our results against identification assumptions.

More generally, estimating IATEs can prove insightful in that some beneficiaries of an agri-environmental policy may exhibit limited or unsatisfactory responses, thereby calling for an intensification of the support, while others may show responses that are well beyond the policy target, thereby suggesting a reduction of the support. Our results illustrate how informative the approach can be in detecting the extent, nature, and source of this heterogeneous response. For instance, we demonstrate that contrasting different farm subgroups can provide additional information on the nature of the heterogeneous response. Specifically, we highlighted that the treatment effect from implementing Pillar 2 agri-environmental measures and fulfilling Pillar 1 conditionality requirements seems more homogeneous than the response to adopting none of the above.

However, the primary policy implication of our results concerns, as mentioned, the need for a better targeting of AEPs. In this respect, caution is necessary as not all farm characteristics considered can

be easily targeted due to practical or political constraints. Nonetheless, our analysis suggests that significant heterogeneity in treatment effects is concentrated within farm sub-groups that can be feasibly targeted. These sub-groups often involve geographical features and specific production specializations. Therefore, delivering some CAP measures at a local scale and tailoring them to specific production orientations, along with broader adoption of results-based payment schemes, may represent a sensible initial step towards better targeting. The new CAP acknowledges greater flexibility for Member States through the new delivery model, allowing them to address the environmental aspects of Pillar 1 (the reinforced CC and the Eco-Schemes replacing the GP) and the AEMs in Pillar 2 more effectively. In principle, this flexibility seems to go along with the goal of improved targeting for these AEPs.

While our empirical results provide valuable insights, our work also contributes to the constructive discussion on the potential and limitations of these relatively new policy assessment methods. In particular, how useful are CML and the analysis of heterogeneous treatment effects in informing policy improvements related to the CAP? Both our conceptual framework and empirical investigation suggest that they can be useful. However, as with all emerging econometric approaches, several issues require careful consideration.

Since standard causal ML methods cannot be used for policy analysis without additional identifying restrictions and assumptions, selecting appropriate confounders and ensuring overlapping/treatment-stable units necessitate a solid theoretical understanding of treatment selection mechanisms. Developing these conceptual foundations also facilitates result interpretation, as the complex output of these estimation methods can be challenging to put into perspective. Among the standard assumptions presented throughout the paper, unconfoundedness and the stable unit treatment value are often regarded as restrictive. Although the former can be corroborated via robustness checks and the use of machine learning algorithms, the latter finds little practical help from the adoption of flexible estimation techniques and thus remains debatable. In this respect, specifying the correct

treatment variable(s) is quintessential for an unbiased interpretation of the resulting treatment effect, an aspect that is still relatively underdiscussed in the literature.

More generally, investigating the effectiveness of CAP's agri-environmental policies within a binary-treatment logic may prove limiting when the analysis targets heterogeneous causal effects. The risk is that the elicited estimates do not entirely reflect farms' heterogeneous responses to a treatment but rather encapsulate the heterogeneity of the treatment itself. Besides the prototypical case of multiple-versions treatments (whether hidden or observable), problems can also arise when a policy measure is not only adopted (i.e., a discrete choice) but also exhibits different intensity levels in different cohorts of farms. In such cases, binary treatments should be extended to incorporate dosage information. How to define the treatment intensity (i.e., the "dose") of different agri-environmental policies is, however, an ambitious empirical question that we leave to future research.

## References

- Afriat, S.N. (1972). Efficiency estimation of production function. *International Economics Review* 13 (3): 568–598.
- Andini, M., Ciani, E., de Blasio, G., D’Ignazio, A., Salvestrini, V. (2018). Targeting with machine learning: An application to a tax rebate program in Italy. *Journal of Economic Behavior & Organization* 156: 86–102.
- Angrist, J. D., Pischke, J. S. (2008). *Mostly Harmless Econometrics*. Princeton University Press, Princeton, NJ.
- Arata, L., Sckokai, P. (2016). The impact of agri-environmental schemes on farm performance in five EU member states: A DID-matching approach. *Land Economics* 92 (1): 167–186.
- Athey, S., Imbens, G.W. (2016). Recursive partitioning for heterogeneous causal effects. *PNAS- Proceedings of the National Academy of Sciences* 113 (27): 7353–7360.
- Athey, S., Imbens, G.W. (2019). Machine learning methods that economists should know about. *Annual Review of Economics* 11: 685–725.
- Athey, S., Tibshirani, J., Wager, S. (2019). Generalised random forests. *Annals of Statistics* 47 (2): 1148–1178.
- Bàrberi, P., Moonen A.C. (2020). Reconciling agricultural production with biodiversity conservation. Burleigh-Dodds Science Publishing, Cambridge (UK).
- Bartolini, F., Vergamini, D., Longhitano, D., Povellato, A. (2021). Do differential payments for agri-environment schemes affect the environment benefits? A case study in Northeastern Italy. *Land Use Policy* 107: 104862.
- Baldoni, E., Coderoni, S., Esposti, R. (2021). Immigrant workforce and agriculture productivity: evidence from Italian farm-level data, *European Review of Agricultural Economics*, 48(4), 805–834, <https://doi.org/10.1093/erae/jbaa033>
- Baldoni, E., Esposti, R., 2020, Agricultural productivity in space: an Econometric Assessment Based on Farm-Level Data. *American Journal of Agricultural Economics* 103(4): 1525–1544.

Belloni, A., Chernozhukov, V., Hansen, C. (2014). High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives* 28 (2): 29–50.

Bertoni, D., Aletti, G., Cavicchioli, D., Micheletti, A., Pretolani, R. (2021). Estimating the CAP greening effect by machine learning techniques: A big data ex-post analysis. *Environmental Science and Policy* 119: 44–53.

Bertoni, D., Curzi, D., Aletti, G., Olper A. (2020). Estimating the effects of agri-environmental measures using difference-in-difference coarsened exact matching. *Food Policy* 90: 101790.

Bodory, H., Busshoff, H., Lechner, M. (2022). High-resolution treatment effects estimation: Uncovering effect heterogeneities with the modified causal forest. *Entropy*, 24(8), 1039.

Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.

Breiman, L. (1984). Classification and regression trees. Routledge. New York. ISBN: 9781315139470.

Brown, C. Kovács, E., Herzon, I., Villamayor-Tomas, S., Albizua, A., Galanaki, A., Grammatikopoulou, I., McCracken, D., Alkan Olsson, J., Zinngrebe, Y. (2021). Simplistic understandings of farmer motivations could undermine the environmental potential of the common agricultural policy. *Land Use Policy* 101, 105136.

Burton, R., Schwarz, G. (2013). Result-oriented agri-environmental schemes in Europe and their potential for promoting behavioural change. *Land Use Policy* 30 (1): 628–641.

Carnegie, N., Dorie, V., Hill, J.J. (2019). Examining treatment effect heterogeneity using BART. *Observational Studies* 5 (2): 52–70.

Carvalho, C., Feller, A., Murraray, J., Woody, S., Yeager, D. (2019). Assessing treatment effect variation in observational studies: results from a data challenge. *Observational Studies* 5 (2): 21–35.

Chabé-Ferret, S., Subervie, J. (2013). How much green for the buck? Estimating additional and windfall effects of French agri-environmental schemes by DID-matching. *Journal of Environmental Economics and Management* 65 (1): 12–27.

Chavas, J.P., Cox, T.L. (1995). On nonparametric supply response analysis. *American Journal of Agricultural Economics* 77 (1): 80–92.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21 (1): C1–C68.

Chipman, H.A., & Gu, H. (2005). Interpretable dimension reduction. *Journal of applied statistics*, 32(9), 969-987.

Chipman, H.A., George, E.I., McCulloch, R.E. (1998). Bayesian CART model search. *Journal of the American Statistical Association* 93 (443): 935–948.

Chipman, H.A., George, E.I., Mcculloch, R.E. (2010). Bart: Bayesian additive regression trees. *Annals of Applied Statistics* 4 (1): 266–298.

Coderoni S., Arata L., Salvatore R., Tiboldo G. (2023). More Enterprise - Youth and female entrepreneurship in agriculture. In: “Mapping and document case studies on family farming in the region of Europe and Central Asia to enhance knowledge exchange through good practices”, (eds.) FAO, Rome (forthcoming).

Coderoni, S., Esposti, R. (2018). CAP payments and agricultural GHG emissions in Italy. A farm-level assessment. *Science of the Total Environment* 627: 427–437.

Coderoni S., Helming J., Pérez-Soba M., Sckokai P., Varacca A. (2021). Key policy questions for ex-ante impact assessment of European agricultural and rural policies. *Environmental Research Letters*, 16 (9): 094044.

Commission of European Communities (2000). Communication from the Commission to the Council and the European Parliament. Indicators for the Integration of Environmental Concerns into the Common Agricultural Policy. COM(2000) 20 final. Commission of the European Communities, Brussels, 26 pp.

Dabkiene, V., Balezentis, T., Streimikiene, D. (2021). Development of agri-environmental footprint indicator using the FADN data: Tracking development of sustainable agricultural development in Eastern Europe. *Sustainable Production and Consumption* 27: 2121–2133.

Davis, J., & Heller, S.B. (2017). Using causal forests to predict treatment heterogeneity: An application to summer jobs. *American Economic Review*, 107 (5): 546–50.

Dessart, F. J., Barreiro-Hurlé, J., van Bavel, R. (2019). Behavioural factors affecting the adoption of sustainable farming practices: a policy-oriented review. *European Review of Agricultural Economics* 46(3): 417–471.

Dick, J., Smith, P., Smith, R., Lolly, A., Moxey, A., Booth, J., Campbell, C., Coulter, D. (2008). Calculating farm scale greenhouse gas emissions. UK Centre for Ecology & Hydrology, Lancaster, UK.

Dorie, V., Harada, M., Carnegie, N. B., & Hill, J. (2016). A flexible, interpretable framework for assessing sensitivity to unmeasured confounding. *Statistics in medicine*, 35(20), 3453-3470.

Dorie, V., Hill, J., Shalit, U., Scott, M., Cervone, D. (2019). Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition. *Statistical Science* 34 (1): 43–68.

Ehlers, M. H., Huber, R., & Finger, R. (2021). Agricultural policy in the era of digitalisation. *Food Policy*, 100, 102019.

Erjavec, K., Erjavec, E. (2015). Greening the CAP – Just a fashionable justification? A discourse analysis of the 2014–2020 CAP reform documents. *Food Policy* 51: 53–62.

Esposti, R. (2000). The impact of public R&D and extension expenditure on Italian agriculture. An application of a mixed parametric/nonparametric approach. *European Review of Agricultural Economics* 27 (3): 365–384.

Esposti, R. (2017a). The empirics of decoupling: Alternative estimation approaches of the farm-level production response. *European Review of Agricultural Economics* 44 (3): 499–537.

Esposti, R. (2017b). The heterogeneous farm-level impact of the 2005 CAP-first pillar reform: A multivalued treatment effect estimation. *Agricultural Economics* 48 (3): 373–386.

Esposti, R. (2022a). The Coevolution of Policy Support and Farmers Behaviour and Performance. An investigation on Italian agriculture over the 2008–2019 period. *Bio-Based and Applied Economics* 11(3): 231–264.

Esposti, R. (2022b). Non-Monetary Motivations of Agri-Environmental Policies Adoption. A Causal Forest Approach. Quaderno di Ricerca n. 459, Dipartimento di Scienze Economiche e Sociali, Università Politecnica delle Marche.

European Commission (2019). The European Green Deal. European Commission, COM(2019) 640 final, Brussels.

European Commission (2020). A Farm to Fork Strategy for a Fair, Healthy and Environmentally-Friendly Food System. COM(2020) 381 final, Brussels.

European Council (2009). Council Regulation (EC) No 1217/2009 of 30 November 2009 setting up a network for the collection of accountancy data on the incomes and business operation of agricultural holdings in the European Community. OJ L 328, 15.12.2009.

Finn, J.A., Bartolini, F., Bourke, D., Kurz, I., Viaggi, D., 2009. Ex-post environmental evaluation of agri-environment schemes using experts' judgements and multicriteria analysis. *Journal of Environmental Planning and Management* 52, 717–737.

Foresight (2011) The Future of Food and Farming Challenges and choices for global sustainability Final Project Report the Government Office for Science, London.

Gelman, A., Stern, H.S., Carlin, J.B., Dunson, D.B., Vehtari, A., Rubin, D.B. (2013). *Bayesian Data Analysis*. Chapman and Hall/CRC, Boca Raton, FL, USA.

Guerrero, S. (2021). Characterising agri-environmental policies: Towards measuring their progress. OECD Food, Agriculture and Fisheries Paper N. 155, OECD Publishing, Paris, France.

Hahn, P.R., Carvalho, C.M., Puelz, D., He, J. (2018). Regularisation and confounding in linear regression for treatment effect estimation. *Bayesian Analysis* 13 (1): 163–182.

Hahn, P.R., Murray, J.S., Carvalho, C.M. (2020). Bayesian regression tree models for causal inference: Regularisation, confounding, and heterogeneous effects (with discussion). *Bayesian Analysis* 15 (3): 965–1056.

Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* Vol. 2, Springer, New York, NY, USA.

Heiler, P., Knaus, M. (2022). Effect or Treatment Heterogeneity? Policy Evaluation with Aggregated and Disaggregated Treatments. IZA Discussion Paper No. 15580, Bonn.

Hill, JL (2011). Bayesian nonparametric modelling for causal inference. *Journal of Computational and Graphical Statistics* 20 (1): 217–240.

Hill, J., Linero, A., Murray, J. (2020). Bayesian additive regression trees: a review and look forward. *Annual Review of Statistics and Its Application* 7: 251–278.

Hill, J., Su, Y.S. (2013). Assessing lack of common support in causal inference using Bayesian nonparametrics: Implications for evaluating the effect of breastfeeding on children’s cognitive outcomes. *The Annals of Applied Statistics*: 1386–1420.

Hoff, P. D. (2009). *A First Course in Bayesian Statistical Methods* Vol. 580, Springer, New York, NY, USA.

Hu, L., Gu, C., Lopez, M., Ji, J., Wisnivesky, J. (2020). Estimation of causal effects of multiple treatments in observational studies with a binary outcome. *Statistical Methods in Medical Research*, 29 (11): 3218–3234.

Huber, W. (2020). Generation of a random variable with fixed covariance structure. Available at: <https://stats.stackexchange.com/questions/444039/whuber-s-generationof-a-random-variable-with-fixed-covariance-structure>

Hudec, B., Kaufmann, C., Landgrebe-Trinkunaite, R., Naumann, S., 2007. Evaluation of Soil Protection Aspects in Certain Programmes of Measures Adopted by Member States. Final report. European Commission, ISBN: 978-92-79-20665-8, doi:10.2779/12302

Hünermund, P., Louw, B., Caspi, I. (2023). Double machine learning and automated confounder selection: A cautionary tale. *Journal of Causal Inference* 11(1): 20220078.

Imai, K., Ratkovic, M. (2013). Estimating treatment effect heterogeneity in randomised program evaluation. *The Annals of Applied Statistics* 7(1): 443–470.

Imbens, G.W. (2020). Potential outcome and directed acyclic graph approaches to causality: Relevance for empirical practice in economics. *Journal of Economic Literature* 58(4): 1129–1179.

Imbens, G.W., Rubin, D.B. (2015). *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press, Cambridge, UK.

Imbens, G.W., Wooldridge, J.M.. (2009). Recent developments in the econometrics of program evaluation. *Journal of Economic Literature* 47(1): 5-86.

IPCC (2006) 2006 Intergovernmental Panel on Climate Change (IPCC) Guidelines for National Greenhouse Gas Inventories prepared by the National Greenhouse Gas Inventories Programme. In: Eggleston H.S., Buendia L., Miwa K., Ngara, T. and Tanabe, K. (eds). Published: IGES, Japan.

Kapelner, A., Bleich, J. (2016). bartMachine: Machine learning with Bayesian additive regression trees. *Journal of Statistical Software* 70 (4): 1–40.

Knaus, M.C., Lechner, M., Strittmatter, A. (2020). Heterogeneous employment effects of job search programmes: A machine learning approach. *Journal of Human Resources*, 0718–9615R1.

Knaus, M.C., Lechner, M., Strittmatter, A. (2021). Machine learning estimation of heterogeneous causal effects: Empirical Monte Carlo evidence. *The Econometrics Journal* 24 (1): 134–161.

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.

Lechner, M. (2018). Modified causal forests for estimating heterogeneous causal effects. *arXiv preprint arXiv:1812.09487*.

Lee, K., Bargagli-Stoffi, F.J., Dominici, F. (2020). Causal rule ensemble: Interpretable inference of heterogeneous treatment effects. *arXiv preprint arXiv:2009.09036*.

Li, F., Ding, P., Mealli, F. (2022). Bayesian causal inference: A critical review. *arXiv preprint arXiv:2206.15460*.

Linden, A., Yarnold, P.R. (2016). Combining machine learning and propensity score weighting to estimate causal effects in multivalued treatments. *Journal of Evaluation in Clinical Practice* 22(6): 871–881.

Lopez, M.J., Gutman, R. (2017). Estimation of causal effects with multiple treatments: A review and new ideas. *Statistical Science* 32 (3): 432–454.

McElreath, R. (2020). *Statistical rethinking: A Bayesian course with examples in R and Stan*. CRC Press.

Mennig, P., Sauer, J. (2020). The impact of agri-environment schemes on farm productivity: A DID-matching approach. *European Review of Agricultural Economics* 47 (3): 1045–1093.

National Rural Network (2010). Report on Cross Compliance Implementation in Italy, Rome.

Nethery, R. C., Mealli, F., & Dominici, F. (2019). Estimating population average causal effects in the presence of non-overlap: The effect of natural gas compressor station exposure on cancer mortality. *The annals of applied statistics*, 13(2), 1242.

Nie, X., Wager, S. (2021). Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika* 108 (2): 299–319.

Ó hUallacháin, D., Finn, J.A., Keogh, B., Fritch, R., Sheridan, H., 2016. A comparison of grassland vegetation from three agri-environment conservation measures. *Irish J. Agric. Food Res.* 55, 176–191.

OECD (2022), *Making Agri-Environmental Payments More Cost Effective*, OECD Publishing, Paris, France.

OECD (2003) *Agriculture and Biodiversity: Developing Indicators for Policy Analysis* Proceedings from an OECD Expert Meeting, Zurich, Switzerland, November 2001.

Pacini, C., Wossink, A., Giesen, G., Vazzana, C., Huirne, R. (2003) Evaluation of sustainability of organic, integrated and conventional farming systems: A farm and field-scale analysis. *Agriculture, Ecosystems & Environment* 95: 273–88.

Paracchini, M.L., Britz, W. (2010) Quantifying effects of changed farm practices on biodiversity in policy impact assessment – An application of CAPRI-Spat OECD, Paris, France.

Pascucci, S., de-Magistris, T., Dries, L., Adinolfi, F., & Capitano, F. (2013). Participation of Italian farmers in rural development policy. *European Review of Agricultural Economics*, 40(4), 605-631.

Pearl, J. (2009). *Causality. Models, Reasoning and Inference*. Cambridge University Press, Cambridge, UK.

Pufahl, A., & Weiss, C. R. (2009). Evaluating the effects of farm programmes: results from propensity score matching. *European Review of Agricultural Economics*, 36(1), 79-101.

Purvis, G., Louwagie, G., Northey, G., Mortimer, S., Park, J., Mauchline, A., Finn, J., Primdahl, J., Vejre, H., Vesterager, J.P., Knickel, K., Kasperczyk, N., Balázs, K., Vlahos, G., Christopoulos, S., Peltola, J. (2009). Conceptual development of a harmonised method for tracking change and evaluating policy in the agri-environment: The Agri-environmental Footprint Index. *Environmental Science and Policy* 12: 321–337.

R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5): 688–701.

Shapley, L. S. (1953). A value for n-person games. *Contributions to the Theory of Games* 2.28 (1953): 307-317.

Stetter, C., Mennig, P., & Sauer, J. (2022). Using Machine Learning to Identify Heterogeneous Impacts of Agri-Environment Schemes in the EU: A Case Study. *European Review of Agricultural Economics* 49(4): 723-739.

Storm, H., Baylis, K., Heckelei, T. (2020). Machine learning in agricultural and applied economics. *European Review of Agricultural Economics* 47 (3): 849–892.

Uehleke, R., Petrick, M., Hüttel, S. (2022). Evaluations of agri-environmental schemes based on observational farm data: The importance of covariate selection. *Land Use Policy* 114, 105950.

VanderWeele, T.J., Hernan, M.A. (2013). Causal inference under multiple versions of treatment. *Journal of Causal Inference* 1(1): 1-20.

VanderWeele, T. J., & Ding, P. (2017). Sensitivity analysis in observational research: introducing the E-value. *Annals of internal medicine*, 167(4), 268-274.

Vanslebrouck, I., Van Huylenbroeck, G., & Verbeke, W. (2002). Determinants of the willingness of Belgian farmers to participate in agri-environmental measures. *Journal of Agricultural Economics*, 53(3), 489–511.

Varacca, A., Arata, L., Castellari, E., & Sckokai, P. (2023). Does CAP greening affect farms' economic and environmental performances? A regression discontinuity design analysis. *European Review of Agricultural Economics*, 50(2), 272-303.

Varian, H. (1984). The nonparametric approach to production analysis. *Econometrica* 52 (3): 579–597.

Vollenweider, X., Di Falco, S., & O'Donoghue, C. (2011). Risk preferences and voluntary agri-environmental schemes: does risk aversion explain the uptake of the Rural Environment Protection Scheme? Grantham Research Institute on Climate Change and the Environment working papers (48). Grantham Research Institute on Climate Change and the Environment, London, UK.

Wager, S., Hastie, T., & Efron, B. (2014). Confidence intervals for random forests: The jackknife and the infinitesimal jackknife. *The Journal of Machine Learning Research*, 15(1), 1625–1651.

Wager, S., Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association* 113 (523): 1228–1242.

Was, A., Malak-Rawlikowska, A., Zavalloni, M., Viaggi, D., Kobus, P., Sulewski, P. (2021). In search of factors determining the participation of farmers in agri-environmental schemes—Does only money matter in Poland? *Land Use Policy*, *101*, 105190.

Wascher, D.M., 2003. Overview of agricultural landscape indicators across OECD countries. Proceedings of the NIJOS/OECD Expert Meeting on Agricultural Landscape, 7–9 October 2002, Oslo, Norway. (Available from URL: <https://research.wur.nl/en/publications/overview-on-agricultural-landscape-indicators-across-oecd-countri>, accessed 25 January 2023).

Westbury, D.B., Park, J.R., Mauchline, A.L., Crane R.T., Mortimer, S.R. (2011). Assessing the environmental performance of English arable and livestock holdings using data from the Farm Accountancy Data Network (FADN). *Journal of Environmental Management* *92* (3): 902–909.

Woody, S., Carvalho, C.M., Murray, J.S. (2021). Model interpretation through lower-dimensional posterior summarisation. *Journal of Computational and Graphical Statistics* *30*(1): 144–161.

Wooldridge, J. M. (2010). *Econometric Analysis of Cross Section and Panel Data*. MIT Press, Cambridge, MA, USA.

Yeager, D.S., Hanselman, P., Walton, G.M., Murray, J.S., Crosnoe, R., Muller, C., Dweck, C.S. (2019). A national experiment reveals where a growth mindset improves achievement. *Nature* *573* (7774): 364–369.

Zhen, H., Qiao, Y., Zhao, H., Ju, X., Zanolli, R., Waqas, M.A., Lun, F., Knudsen, M.T. (2022). Developing a conceptual model to quantify eco-compensation based on environmental and economic cost-benefit analysis for promoting ecologically intensified agriculture. *Ecosystem Services* *56*, 101442.

Zhou, T., Elliott, M. R., & Little, R. J. (2019). Penalized spline of propensity methods for treatment comparison. *Journal of the American Statistical Association*, *114*(525), 1-19.

Zimmermann, A., Britz, W. (2016). European farms' participation in agri-environmental measures. *Land Use Policy* *50*: 214–228.

**Table 1.** The policy treatments set

<b>Treatment</b>	<b>Agri-environmental policies</b>	<b>Total number of farms (2015–2018 FADN balanced sample)</b>
$T_1$	None	512
<i>Control Group</i>	Only those implied by the Pillar 1 DPs	2,841
$T_2$	Those implied by both Pillar 1 DPs and Pillar 2 AEMs.	648

**Table 2.** Elementary indicators used to assemble the outcome variable (FADN-AFI)

<b>Environmental issue</b>	<b>Assessment criterion</b>	<b>Indicator</b>	<b>Measurement Unit</b>	<b>Weight</b>
Natural resources protection	Intensity of crop husbandry/livestock production	Fertiliser cost per ha of utilised agricultural area (UAA)	Euro/ha	-1
		Crop protection costs per ha UAA	Euro/ha	-1
		Average number of livestock units per hectare	LU/ha	-1
		Energy consumption	Energy costs per hectare UAA	Euro/ha
		% of UAA irrigated	Share	-1
Biodiversity & land use	Land use diversity	Crop diversity (Shannon diversity) index ( $BI_{it}$ )	Index (0–1)	+1
		Provision of woodland habitats	% of total farm area that is woodland	Share
Climate change mitigation	GHG emissions	GHG at the farm level	kg CO <sub>2eq</sub>	-1

**Table 3.** List of covariates used in the analysis

	Unit of Measure	Component	Reference
<i>Economic characteristics</i>			
Total Arable Land	hectares	$V_{it}$	(1), (2), (3), (4), (6), (7), (8)
Share of Rented Land	%	$V_{it}$	(2), (3), (4), (7), (8)
Farm Revenue	€ per hectare	$V_{it}$	(4), (7), (8)
Farm Fixed Costs	€ per hectare	$V_{it}$	New
Farm Variable Costs	€ per hectare	$V_{it}$	(4)
Fertilizer Expenditure	€ per hectare	$V_{it}$	(2), (4), (7)
Pesticides Expenditure	€ per hectare	$V_{it}$	(2), (4), (7)
Livestock Density	Units per hectare	$V_{it}$	(2), (5), (6), (8)
Family Labour	count	$S_i^{(*)}$	(2), (3), (4), (8), (9)
Non-Family labour	count	$V_{it}$	(2), (3), (4), (9)
Share of the Most Important Crop	%	$V_{it}$	(5), (6)
Share of the Second Most Important Crop	%	$V_{it}$	(5), (6)
Share of Grassland	%	$V_{it}$	(4), (7), (8)
Machinery Horsepower	Kw per hectare	$V_{it}$	(3)
Machinery Value	€ per hectare	$V_{it}$	(3)
Machinery Endowment	Units per hectare	$V_{it}$	(3)
Farm Specialisation	categorical	$S_i$	(3), (6), (7), (8)
<i>Socio-demographic characteristics</i>			
Farmer's Age	years	$S_i^{(*)}$	(1), (3), (5), (6), (7), (8), (9)
Farmer's Gender	categorical	$S_i$	New
Farmer's Education	categorical	$S_i^{(*)}$	(1), (3), (8)
Experience with Previous AEPs	categorical	$S_i$	(6), (7)
<i>Environmental/geographical characteristics</i>			
Disadvantaged Area	categorical	$S_i$	(5), (7), (8)
Latitude and Longitude	degrees	$S_i$	(6), (7)
Average Altitude	meters	$S_i$	(9)

(\*) Although these characteristics are not strictly time invariant, we assume they are approximately so ( $V_{it} \approx S_i$ ) for the period 2015–2018.

The reference papers are: (1) Vanslebrouck et al. (2002); (2) Pufhal and Weiss (2009); (3) Pasucci et al. (2013); (4) Arata and Sckokai (2016); (5) Zimmerman and Britz (2016); (6) Bertoni et al. (2020); (7) Uehleke et al. (2022); (8) Wąs et al., 2021; (9) Varacca et al. (2023).

**Table 4.** IATEs estimates for Model (6).  $0 \notin \text{CrI}$  indicates the proportion of IATEs' CrI that do not include zero.  $\text{MAP} > (<) 0$  denotes the proportion of IATE with MAP greater (lower) than zero.  $\text{ATE} > (<) 0$  indicates the posterior probability that the ATE (as defined in Eq. (4)) is greater (lower) than zero. CrI ATE is the 95% CrI for the ATE.

Treatment	$0 \notin \text{CrI}$	$\text{MAP} > 0$	$\text{MAP} < 0$	$\text{ATE} > 0$	$\text{ATE} < 0$	ATE	CrI ATE	
							lower	upper
<i>All observations</i>								
	(%)	(%)	(%)	(%)	(%)	(AFI)	(AFI)	(AFI)
$T_2$	72.3	100	0	100	0	0.55	0.31	0.79
$T_1$	15.3	28.5	71.5	0.5	99.5	-0.42	-0.73	-0.08
<i>Observations for which <math>0 \notin \text{CrI}</math></i>								
$T_2$	1	100	0	100	0	0.61	0.364	0.86
$T_1$	1	0	100	0	100	-0.85	-1.02	-0.46

<sup>i</sup> Therefore, we here consider AEMs as a subset of whole menu of AEPs.

<sup>ii</sup> Measure 10 supports, among others: integrated production, manure management, increasing soil organic matter, sustainable management of extensive grassland, and management of buffer strips against nitrates. Measure 11 supports both conversion to and maintenance of organic practices and methods. It is worth noticing that Stetter et al. (2022: 732) do not consider the organic farming measure “due to [the] distinctly different farming approach compared to conventional farms”. As clarified in Section 4, here we include this measure in the analysis in order to perform a comparison with the results obtained on the whole sample.

<sup>iii</sup> At the Member-State level, the total amount of GP must correspond to 30% of the total DPs. In several EU countries (including Italy), this condition is satisfied by automatically assigning to eligible farms 30% of total DP as the GP.

<sup>iv</sup> Since production decisions must be taken ex-ante, their consequences are evidently subject to some degree of uncertainty. Consequently, farmers actually maximize  $E\{g(T_{it,k}, \mathbf{X}_i)\}$  and, more importantly, the condition  $E\{\Pi[g(T_{it,k}, \mathbf{X}_i)]\} \geq E\{\Pi[g(T_{it,h}, \mathbf{X}_i)]\}, \forall k, h \in K, k \neq h$  remains valid only if we are willing to assume farmer’s risk neutrality. Otherwise, the variance of  $\pi_{it,k}$  and  $\pi_{it,h}$ , and the possible impact of  $T_{it,k}$  on them, would also matter.

<sup>v</sup> It can be argued that, under risk aversion, farmers are expected to be more prudent and conservative; therefore, *ceteris paribus*, the participation in the treatment and the observed response,  $\Delta \mathbf{y}$ , should be smaller. At the same time, however, the monetary support granted to participant farmers may represent a guaranteed income, making participation in the measure a less risky situation. Also notice that, under risk aversion, risk itself can be interpreted as an additional source of costs and/or foregone income that the AEP is expected to compensate. Therefore, as noted in previous studies (Esposti 2017a,b), it is difficult to model and predict the differential impact of these support measures between risk-neutral and risk-averse farmers.

<sup>vi</sup> Unlike the other vectors of model variables, the netput vector is here indicated with a small letter,  $\mathbf{y}_{it}$ , to avoid confusion with the conventional notation of potential outcomes, i.e.  $Y_i(0)$  and  $Y_i(1)$  (see Section 5).

<sup>vii</sup> As will be clarified in Section 4.4, examples of internal factors are the farm size and the farmer’s age and education. Examples of external factors are latitude and farm’s location in a disadvantaged area.

<sup>viii</sup> Following the conventional terminology of production theory, this should be a direct profit function as opposed to the more frequently used indirect profit function, where profit is a function of only output and input prices. In fact, in addition to netput quantities, the direct profit function includes the respective prices expressed as  $\Pi[\mathbf{v}'_{it}g(T_{it,k}, \mathbf{X}_i)]$  where  $\mathbf{v}'_{it}$  is the  $(M \times 1)$  vector of netput prices. For non-market netputs, there are no prices but these elements in  $\mathbf{v}'_{it}$  can still be interpreted as shadow prices. Nonetheless, prices have been excluded from the present notation under the assumption, maintained that the prices are constant or, more precisely, unaffected by the policy regime.

<sup>ix</sup> The heterogeneity among farms is the core of this theoretical framework. With homogeneous farms, we would have  $\pi_{it,k} = \pi_{jt,k} = \pi_t, \forall i \neq j, \forall k$  and  $\forall t$ , so all farmers would opt for the same policy, and we would observe only one treatment. A policy response would thus be only conjectural but not actually observable if not by comparing farms before and after the treatment.

<sup>x</sup> Notice that this assessment applies to both single treatment and multiple treatments versions.

<sup>xi</sup> See Zimmerman and Britz (2016), Dessart et al. (2019), Brown et al. (2021) for recent and extensive reviews of both structural and behavioural factors underlying farmer’s decisions.

<sup>xii</sup> The programming period has been subsequently extended to 2022, also as a consequence of the COVID-2019 pandemic. Validated data on years 2021 and 2022 have still to be released.

<sup>xiii</sup> It is worth noticing that, extracting the balanced sample from the unbalanced one does not imply a relevant loss in terms of representativeness of the sample (see Baldoni et al., 2021 for a detailed explanation).

<sup>xiv</sup> More specifically, from a survey carried out at national level, it emerged that there are 65 different versions of the Measure 10, that can be applied at regional programming level, corresponding to a total of 100 commitment categories for the whole 21 RDPs (Source: <https://www.reterurale.it/flex/cm/pages/ServeBLOB.php/L/IT/IDPagina/23816> accessed on 05/06/2023).

<sup>xv</sup> The support for organic farming is exemplary in this respect. The nature of the response may vary largely across different farming types, even under such a very specific measure. The same argument applies to CC requirements, where each element and constraint becomes applicable to the farm depending on the characteristics of the farmland or the agricultural activities carried out.

<sup>xvi</sup> For elements of  $\mathbf{y}_{it,k}$  that are only marginally (or not at all) affected by the policy treatment under consideration, we have  $\Delta \mathbf{y}_{it,k} \approx 0$ . Therefore, we may restrict the analysis only to input and output decisions that are related to the environmental measures, all the rest being orthogonal by assumption.

<sup>xvii</sup> These goals are related to: (i) the mandatory practices devised to benefit the environment (soil and biodiversity in particular) and climate (with the GP of Pillar 1); (ii) the new RDPs’ priority areas specifically addressing the environment and climate change (Pillar 2). The latter are aimed at ‘Restoring, preserving and enhancing ecosystems dependent on agriculture and forestry’ (Priority 4) and ‘Promoting resource efficiency and supporting the shift toward a low-carbon and climate-resilient economy in the agriculture, food and forestry sectors’ (Priority 5).

<sup>xviii</sup> See Appendix Table B2 for further details.

<sup>xix</sup> Following Purvis et al. (2009), all the indicators and assessment criteria in the FADN-AFI receive a subjectively equal weighting.

<sup>xx</sup> Averaging only over the last two years reduces the risk of integrating out potential ‘accumulation effects’ by smoothing over a longer period (i.e., the cumulative benefit of environmentally friendly practices).

<sup>xxi</sup> This explains the presence of ‘Insurance Expenditure’ among covariates. This variable might seem contradictory to the risk neutrality assumed in deriving the theoretical framework (Section 3). However, it is worth remembering that, in most cases, farms incur these

---

costs not because of their risk aversion but because taking out an insurance contract is mandatory in order to receive public or private investment support. For this reason, this variable was considered in previous studies and, consequently, also in the present study.

<sup>xxii</sup> This requires assuming no anticipation and no instantaneous impact of either  $T_1$  or  $T_2$  on  $\mathbf{V}_{it}$ . With no anticipation, we refer to the assumption that farmers have not changed their characteristics  $\mathbf{V}_{it-1}$  in response to the foreseen implementation of the policy at time  $t$ .

<sup>xxiii</sup> In short, the authors discuss how a construct resulting from the interaction between farm type, farm size, farmer's age, farm capital intensity and proxies for risk behaviour is conceivably strongly correlated with the unobservable trait, thereby contributing to de-confounding the treatment effect.

<sup>xxiv</sup> For an inventory of these methods, see Nie and Wager (2021).

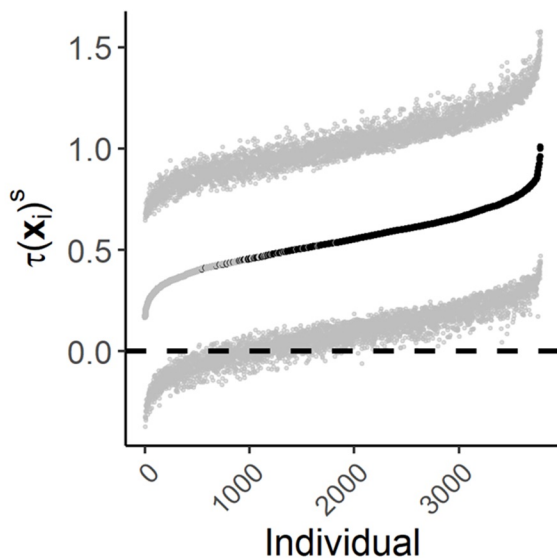
<sup>xxv</sup> Notice that, although the terminology 'Casual Forests' resembles that used in Wager and Athey (2018), BCF differ substantially from the frequentist counterpart in their definition, functioning, and in how inference is performed.

<sup>xxvi</sup>  $f(\mathbf{X}_i, T_i)$  could be specified as a fully parametric function, although this would inevitably constraint the cross-farm technological and behavioural heterogeneity. Admitting an arbitrarily complex function is thus more consistent with the assumption of a farm-specific production set  $F_i$ .

<sup>xxvii</sup> The importance metric is actually obtained from a BART which includes the PS (the so-called PS-BART). Unlike the algorithm in equation (7), the PS-BART does not distinguish the prognostic from the treatment effect component. However, in terms of variable importance, the difference between the two techniques is negligible.

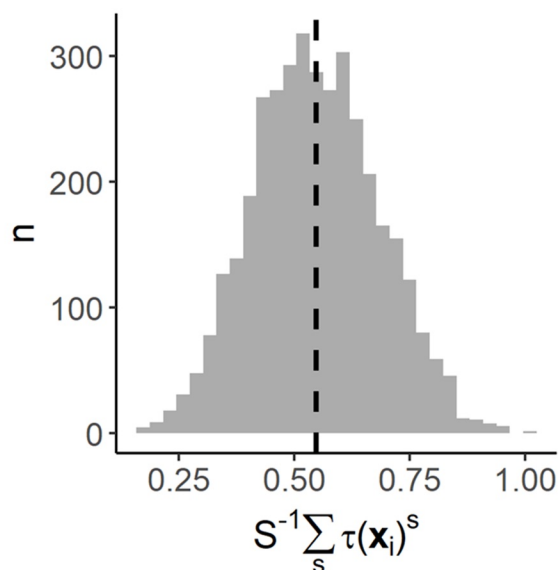
## A Treatment group $T_2$ , i.e.: adopting AEM.

**1** IATEs (MAP + 95% HPDIs)



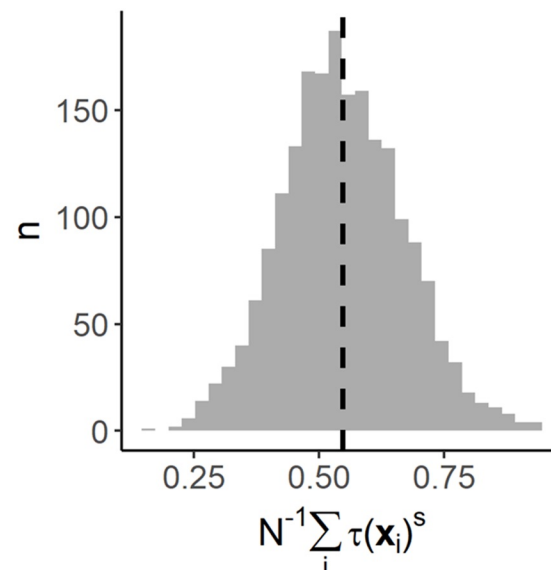
**2** Distribution of IATEs (MAP)

$$N^{-1}S^{-1}\sum_s\sum_i\tau(\mathbf{x}_i)^s = 0.55$$



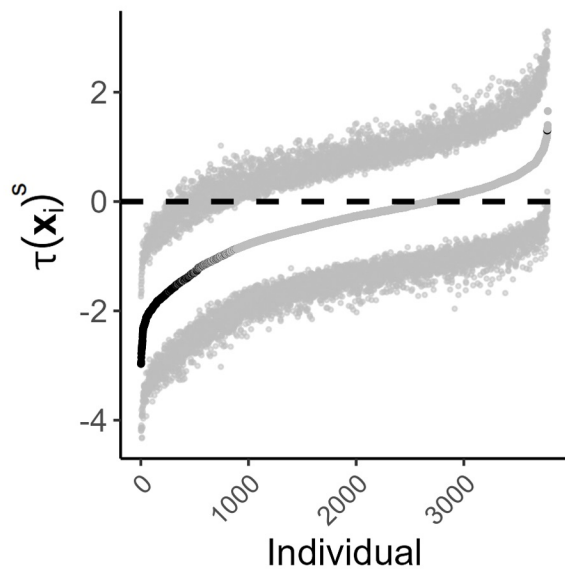
**3** Posterior Distribution of ATE

$$N^{-1}S^{-1}\sum_s\sum_i\tau(\mathbf{x}_i)^s = 0.55$$



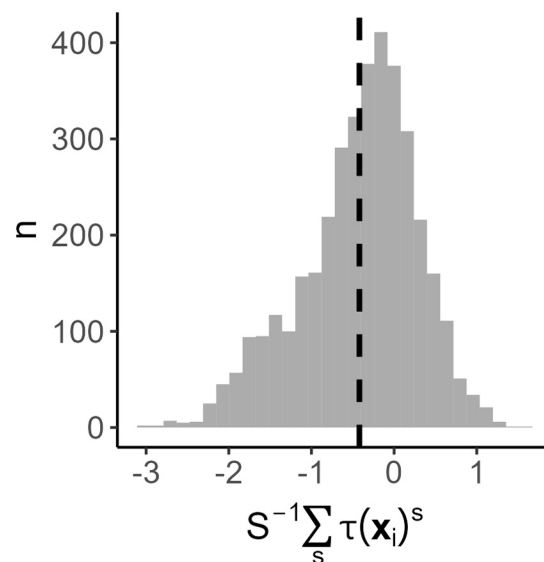
## B Treatment group $T_1$ , i.e.: farmers not fulfilling any conditionality constraints nor implementing AEM.

**1** IATEs (MAP + 95% HPDIs)



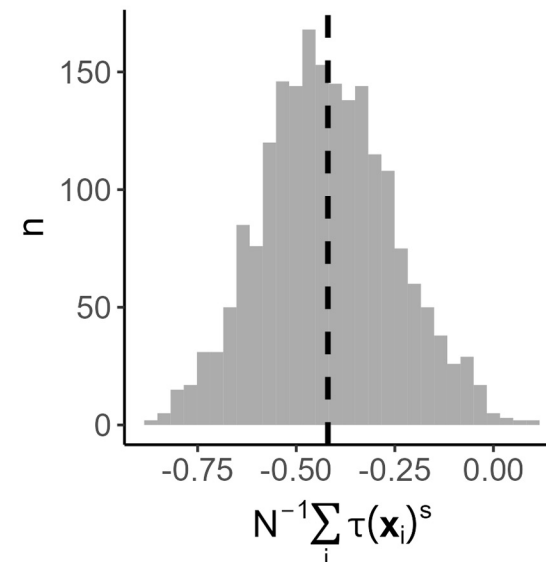
**2** Distribution of IATEs (MAP)

$$N^{-1}S^{-1}\sum_s\sum_i\tau(\mathbf{x}_i)^s = -0.42$$



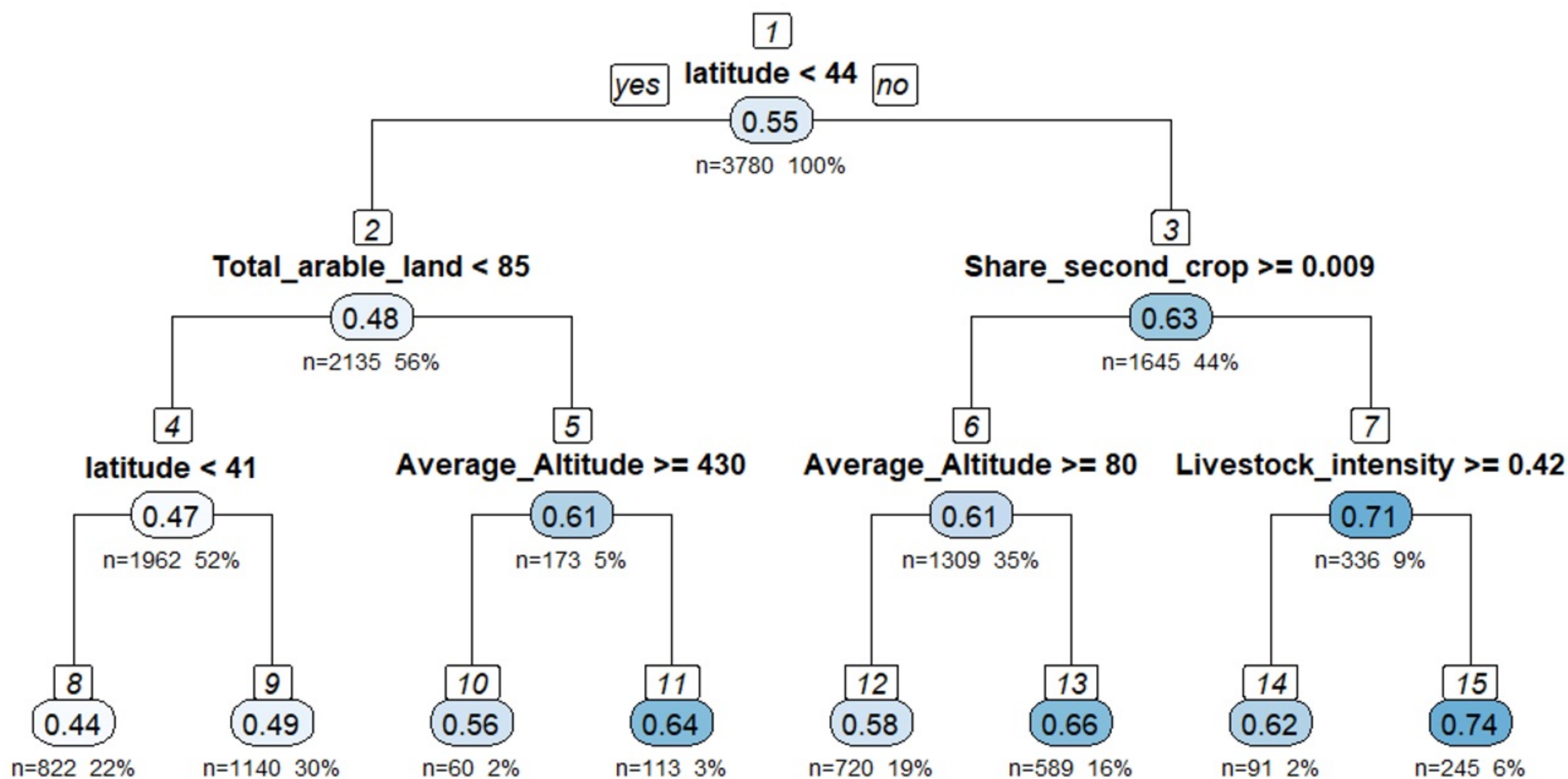
**3** Posterior Distribution of ATE

$$N^{-1}S^{-1}\sum_s\sum_i\tau(\mathbf{x}_i)^s = -0.42$$

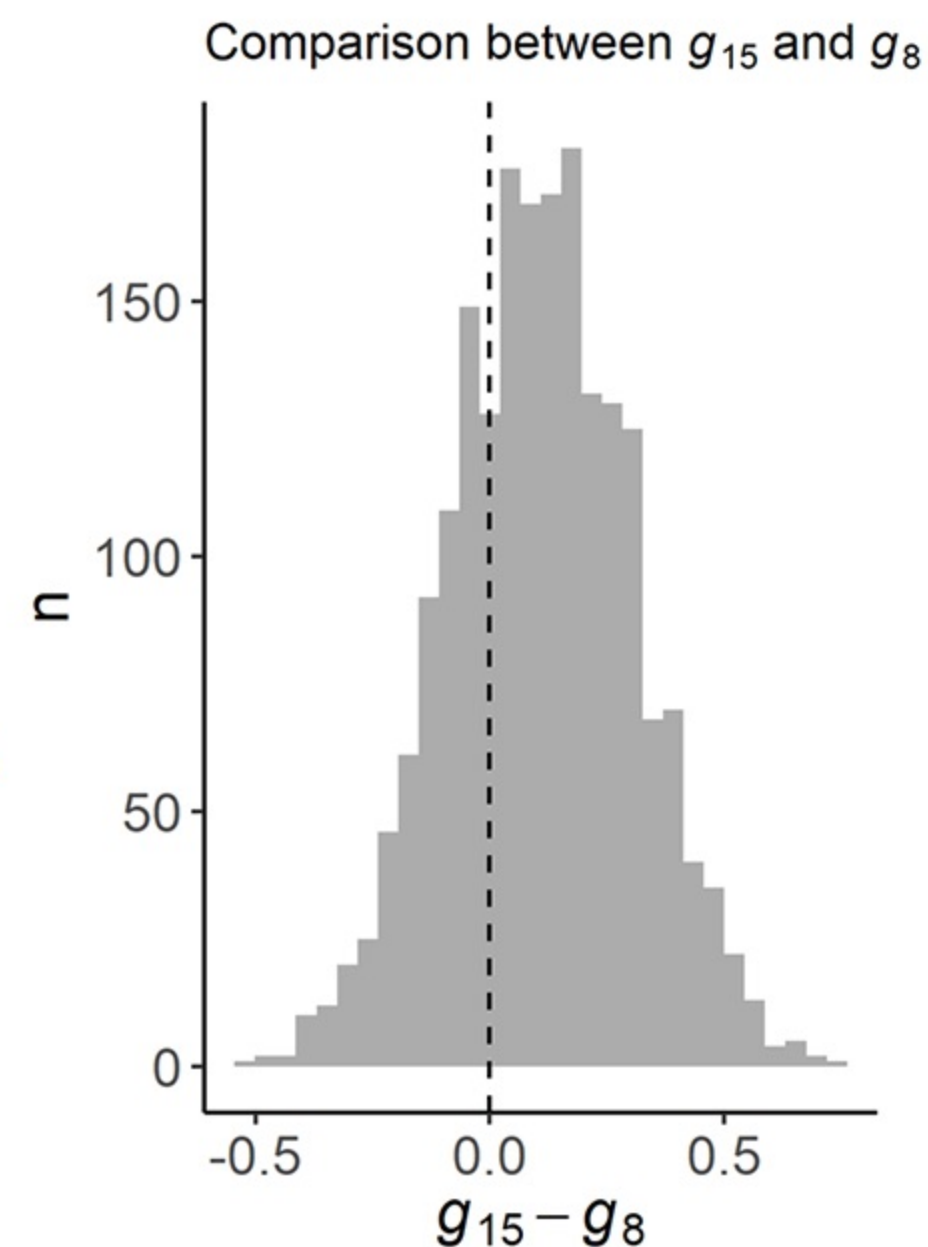


**A** Treatment  $T_2$ . I: structure of the penalised regression tree. II: posterior distribution for the comparison between subgroup 15 and subgroup 8.

1

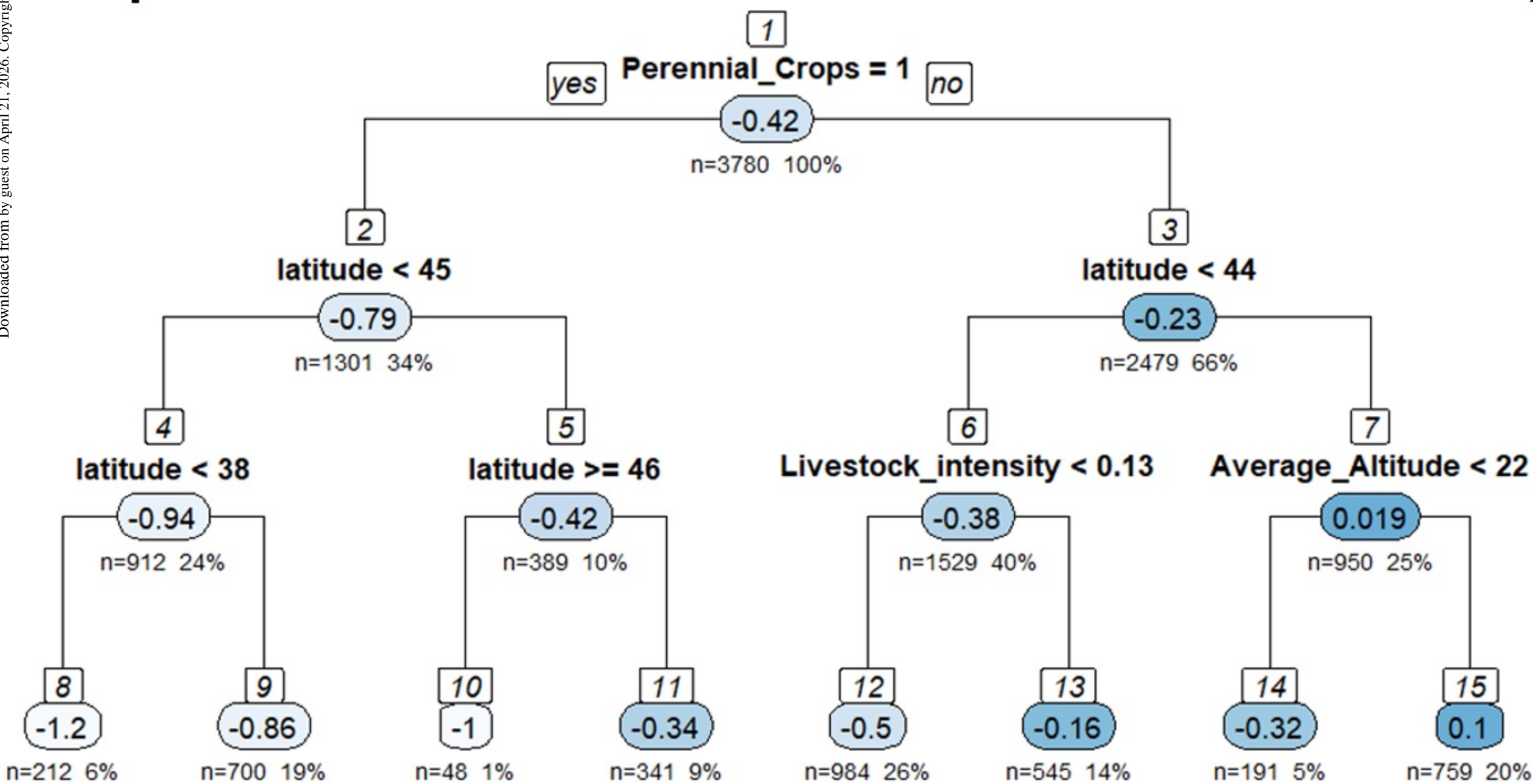


2

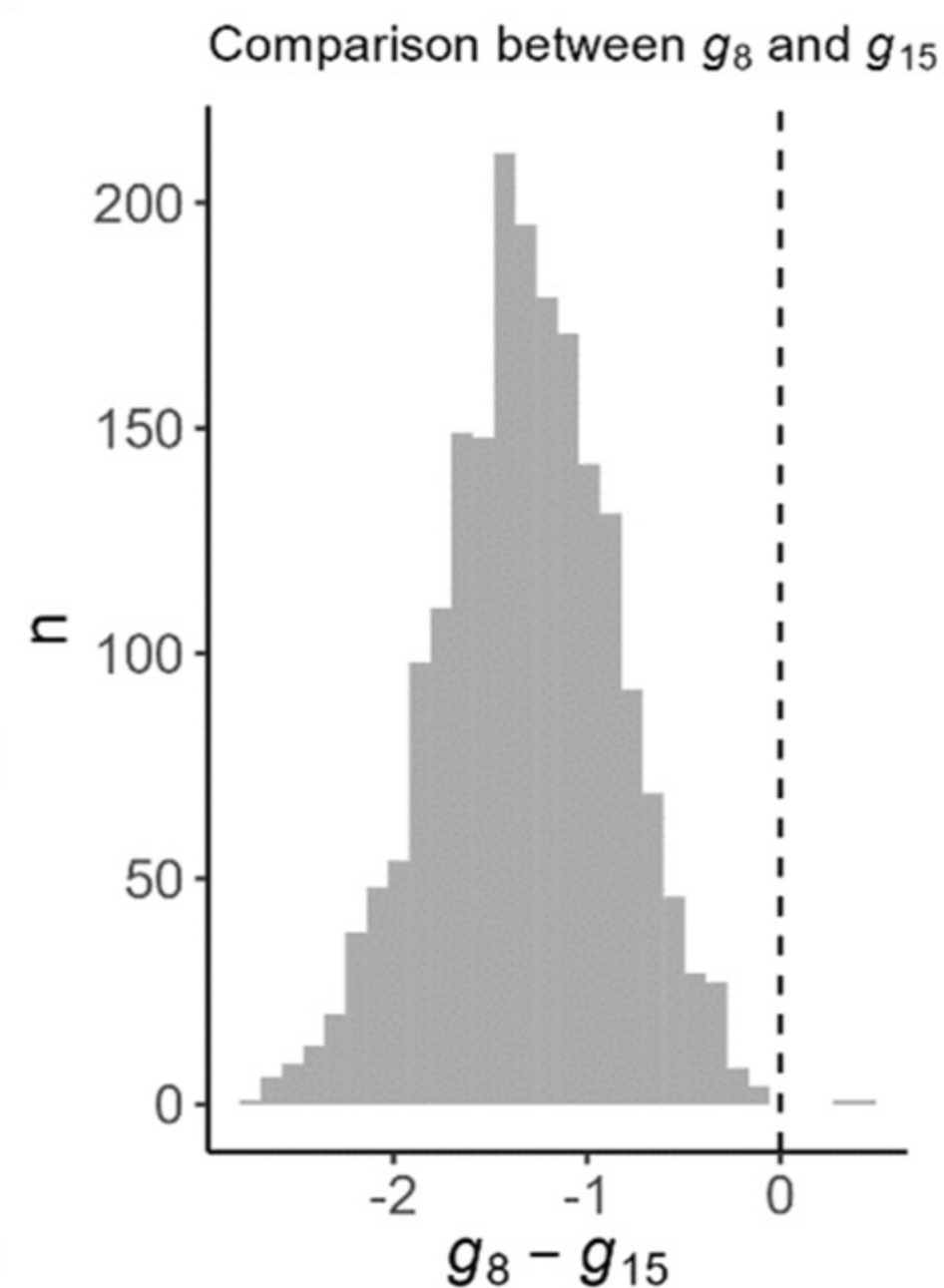


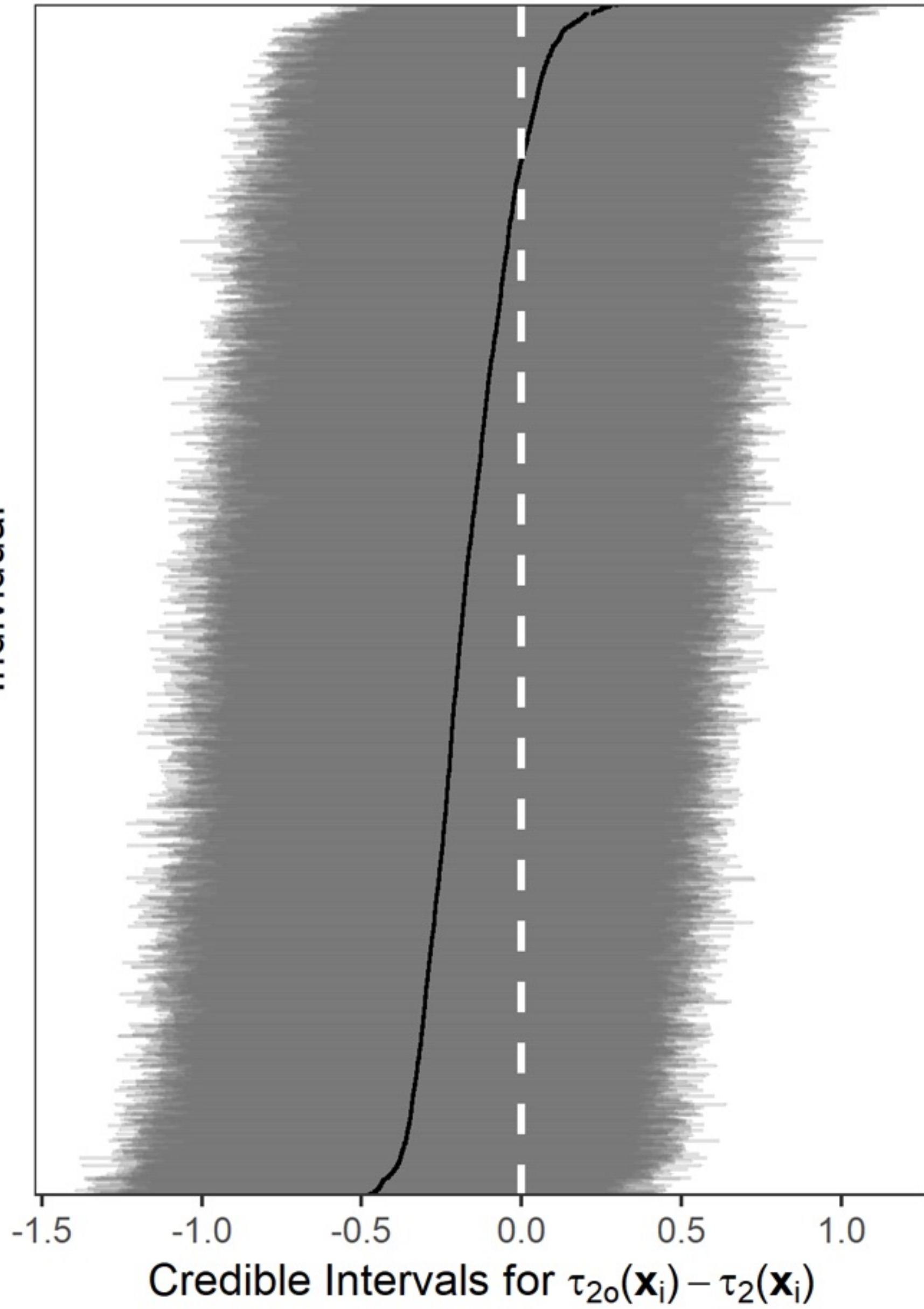
**B** Treatment  $T_1$ . I: structure of the penalised regression tree. II: posterior distribution for the comparison between subgroup 15 and subgroup 8.

1



2



**1****2**