

Generalized Causal Tree for Uplift Modeling*

Preetam Nandy^{1†}, Xiufan Yu^{2†}, Wanjun Liu³, Ye Tu⁴, Kinjal Basu⁴, Shaunak Chatterjee⁴

¹LinkedIn Corporation, ²University of Notre Dame, ³LinkedIn Corporation, ⁴Aliveo AI Corp

Abstract

Uplift modeling is crucial in various applications ranging from marketing and policy-making to personalized recommendations. The main objective is to learn optimal treatment allocations for a heterogeneous population. A primary line of existing work modifies the loss function of the decision tree algorithm to identify cohorts with heterogeneous treatment effects. Another line of work estimates the individual treatment effects separately for the treatment group and the control group using off-the-shelf supervised learning algorithms. The former approach that directly models the heterogeneous treatment effect is known to outperform the latter in practice. However, the existing tree-based methods are mostly limited to a single treatment and a single control use case, except for a handful of extensions to multiple discrete treatments. In this paper, we propose a generalization of tree-based approaches to tackle multiple discrete and continuous-valued treatments. We focus on a generalization of the well-known causal tree algorithm due to its desirable statistical properties, but our generalization technique can be applied to other tree-based approaches as well. The efficacy of our proposed method is demonstrated using experiments and real data examples.

Keywords: Continuous treatment; Heterogeneous causal effect; Multiple treatments; Randomized experiments; Targeted advertising; Personalized recommendations.

*This work was initiated when P. Nandy, Y. Tu, K. Basu, and S. Chatterjee were at LinkedIn, and X. Yu was interning at LinkedIn. P. Nandy is now at Google.

[†]The first two authors contributed equally.

1 Introduction

Uplift modeling (Gutierrez and Gérardy, 2017; Zhang et al., 2021) aims to identify individuals who benefit the most as a result of receiving a certain intervention. It has a wide range of applications, including political or marketing campaigns, policy-making, ad targeting, personalized medicine, and much more. Under Rubin’s framework of causal inference (Rubin, 1974), uplift modeling amounts to estimating the Conditional Average Treatment Effect (a.k.a. heterogeneous causal effect or individual treatment effect) for each individual based on their characteristics (or heterogeneity factors). Uplift modeling usually considers randomized control experiment data with a treatment group and a control group. The treatment group is a randomly selected population segment where each individual receives the intervention (or treatment) under consideration (e.g., medical treatment or a marketing campaign). The control group is another disjoint random segment to which the intervention is not applied. The aim is to predict the difference in the counterfactual responses when an individual receives treatment versus when the same individual receives control. These estimates are then used to decide whether the treatment should be applied. An important modification to the problem is when several treatments (possibly infinitely many) are available to us. In this case, we need to decide whether a treatment should be applied or not and which treatment (e.g., medicine dose) is the most beneficial in a given case.

Though there is a huge literature on uplift modeling, most of the focus has been on binary treatments, that is, one treatment and one control; see Zhang et al. (2021) for a recent literature review. It is of utmost importance to develop methods for uplift modeling with multiple discrete and continuous treatments, e.g., determining continuous weight parameters for search or recommender systems (Agarwal et al., 2018), choosing different doses of a drug (Jin and Rubin, 2008), choosing various thresholds for sending notifications or emails (Gupta et al., 2016; Liu et al., 2023). Recently, there are some methods that have tried to tackle the finitely many multiple-treatment (Rzepakowski and Jaroszewicz, 2012; Zhao et al., 2017;

Guo et al., 2020; Zhou et al., 2023) or continuous-treatment (Oprescu et al., 2019; Bica et al., 2020a; Wan et al., 2022) scenarios. In this paper, we propose a novel generalization of the well-known Causal Tree (Athey and Imbens, 2016) algorithm, called Generalized Causal Tree (GCT), for uplift modeling with multiple discrete or continuous treatments. We chose to generalize the causal tree algorithm due to its desirable statistical properties, but our generalization technique can also be applied to other decision-tree-based approaches.

The paper is organized as follows. Section 2 introduces the problem setup and the causal tree approach. Section 3 discusses the details of the GCT algorithm and an implementation-friendly variant of GCT. The output of GCT is user-friendly in the sense that it is (i) interpretable for understanding the treatment effect diversity and (ii) easily utilizable for optimal treatment allocations. We demonstrate the efficacy of our algorithm through empirical evaluations in Section 4 and uplift analysis on two real-world datasets in Section 5 before concluding with a discussion in Section 6. We end this section by highlighting some of the related works.

Related Works: The uplift modeling problem can be decomposed into several sub-regression problems that can be solved with any supervised learning method. These types of approaches are called meta-algorithms or meta-learners (Künzel et al., 2019). A sub-class of meta-learners are the so-called two-model approaches (a.k.a. T-learner) that build two separate predictive models for the treatment and control groups. The two-model approach naturally extends to continuous treatments (or multiple discrete treatments) by modeling the treatment group response as a function of the covariates and the treatment value. The treatment indicator plays a different role than the other features in the two-model approach, but there exist meta-learners where the treatment indicator is considered as an additional covariate (or feature), and the treatment effect predictions are computed from a single supervised learning model (Hill, 2011; Green and Kern, 2012). These single-model based meta-learners are also known as S-learners, and they can also be extended to continuous treatments. There are other meta-learner approaches that use special formulations for binary treatments, in-

cluding X-learner (Künzel et al., 2019) and PW-learner (Curth and van der Schaar, 2021). The main advantage of this approach is that state-of-the-art machine learning models (such as Random Forest and XGBoost) can be used directly for the “uplift” prediction. However, good prediction performances for both the treatment and control group models do not guarantee good uplift predictions. One of the main reasons is that the most relevant variables for the treatment model or the control model might not be the most relevant variables for the uplift predictions. We refer to Section 5 of Radcliffe and Surry (2011) for a simulation-based illustration along with a nicely presented list of intuitive arguments on the failure of the two-model approach. For further empirical evidence on the disadvantage of the two-model approach, see Knaus et al. (2020); Zaniewicz and Jaroszewicz (2013).

Another line of work models the treatment effect directly by modifying well-known classification machine learning algorithms. For example, Guelman et al. (2014) focused on k -nearest neighbors while Zaniewicz and Jaroszewicz (2013) proposed a modification of the SVM model. Other methods in the literature are based on modifications of the splitting criterion in decision trees (Athey and Imbens, 2016; Radcliffe and Surry, 2011; Rzepakowski and Jaroszewicz, 2012). Hansotia and Rukstales (2002) proposed a splitting criterion based on the difference in the uplift between two leaves. Radcliffe and Surry (2011) proposed to fit a linear model with an interaction term representing the relationship between treatment and the split and defined a splitting criterion based on the significance of the interaction term. Rzepakowski and Jaroszewicz (2012) described three splitting criteria based on the divergence between the treatment and control groups in the leaves. Athey and Imbens (2016) proposed the causal tree algorithm through a treatment effect based splitting criterion. Most of these approaches are focused on binary uplift models, i.e., single treatment and single control use case. A natural extension to multiple treatment cases is to build binary uplift models for each treatment separately. However, this simple extension is infeasible for continuous treatments. Moreover, building separate models for each treatment is computationally and statistically inefficient. To address these issues, some authors (Rzepakowski

and Jaroszewicz, 2012; Zhao et al., 2017) proposed appropriate adjustments to the splitting criteria that can tackle multiple treatments simultaneously in a modified decision tree approach. Unfortunately, these extensions cannot handle continuous treatments without an ad hoc discretization of treatment values. Our main contribution is to show that there is no need for an ad hoc discretization of continuous treatment values when we can learn an optimal segmentation of the treatment values from data.

2 Uplift Modeling using Causal Tree

We define uplift modeling as a problem of optimal treatment selection to maximize future average response in a population using randomized experiment data. We discuss how the causal tree algorithm can be used to maximize the average response in the single treatment case. Next, we point out the main disadvantages of the existing causal-tree related approaches for the multiple discrete and continuous treatment use cases.

Problem Definition: We consider a randomized experiment with control and treatment groups. The data consists of the treatment assignment $T \in \{0\} \cup \mathcal{Z}$, the response variable $Y \in \mathbb{R}$, and potential heterogeneity factors $\mathbf{X} \in \mathcal{X}$, where \mathcal{Z} and \mathcal{X} denote the set of treatment values in the treatment group and the set of all values taken by \mathbf{X} respectively. Without loss of generality, we assume that each member in the control group received $T = 0$, and each member in the treatment group is assigned a treatment value $Z \in \mathcal{Z} \subseteq \mathbb{R} \setminus \{0\}$ generated from a (discrete or continuous) probability distribution $p_Z(\cdot)$ on \mathcal{Z} . For the single treatment case, $p_Z(\cdot)$ is a point mass distribution on $\mathcal{Z} = \{1\}$ (without loss of generality). In this paper, we do not discuss the choice of \mathbf{X} but focus on exploiting the heterogeneity of the conditional treatment effect to maximize the average response for a given \mathbf{X} . The set of potential heterogeneity factors \mathbf{X} are assumed to be pre-treatment covariates (i.e., \mathbf{X} is independent of T).

Let $\mu(t, \mathbf{x}) = E[Y(T = t) \mid X = \mathbf{x}]$ denote the average counterfactual response for individuals with characteristics (or heterogeneity factors) $\mathbf{X} = \mathbf{x}$ under the treatment $T = t$.

We further define $\tau(z, \mathbf{x}) = \mu(z, \mathbf{x}) - \mu(0, \mathbf{x})$ to be the conditional treatment effect for $\mathbf{X} = \mathbf{x}$ and $z \in \mathcal{Z}$. For single treatment case, we denote the individual treatment effect by $\tau(\mathbf{x})$.

Causal Tree: The single treatment causal tree algorithm of Athey and Imbens (2016) fits a decision tree model for estimating conditional treatment effect $\tau(\mathbf{x})$ with a splitting criterion that maximizes the sum of squared treatment effects in the leaves (i.e., nodes without children) while penalizing for higher estimation variances. The output of the causal tree algorithm is a partition (i.e., a disjoint and exhaustive collection) Π_1, \dots, Π_K of \mathcal{X} and the estimated causal effects $\hat{\tau}(\Pi_1), \dots, \hat{\tau}(\Pi_K)$ in those cohorts. Based on the estimated causal effects, one can choose treatment or control for each cohort for future assignments, where the within cohort treatment effects are assumed to be homogeneous (or statistically indistinguishable). Tu et al. (2021) used the single treatment causal tree algorithm in a multi-treatment uplift modeling setting by building separate trees for each treatment and then merging the trees to obtain a finer partition Π'_1, \dots, Π'_K of \mathcal{X} . The authors also described a procedure to obtain cohort-specific treatment effect estimates $\tau(z, \Pi'_i)$ for all $z \in \mathcal{Z}$ using the estimated treatment effects in the separately built trees. These estimates can be used to maximize the future total effect by assigning

$$v(\mathbf{x}) = \begin{cases} \arg \max_{z \in \mathcal{Z}} \hat{\tau}(z, \Pi'_i) & \text{if } \max_{z \in \mathcal{Z}} \hat{\tau}(z, \Pi'_i) > 0, \\ 0 & \text{otherwise} \end{cases}$$

to each Π'_i for $i = 1, \dots, K$.

Sample Inefficiency: This approach can suffer from sample inefficiency by ignoring the presence of similarities between treatments. For example, suppose treatments $T = 1$ and $T = 2$ are identical, then there is no need to learn separate causal trees for $T = 1$ and $T = 2$ by splitting the data corresponding to $T \in \{1, 2\}$. Of course, we may not have prior knowledge of treatment similarities. Still, we can be more sample efficient by learning treatment cohorts (i.e., a similarity-based partition of the treatment set \mathcal{Z}) jointly with the \mathbf{X} cohorts. Moreover, the partitioning of the treatment set is a necessity when \mathcal{Z} contains

infinitely many values (e.g., continuous treatment), as building separate trees is not possible anymore.

Continuous Treatment: Existing continuous treatment version of the causal tree or causal forest (Wager and Athey, 2018) algorithm defines the treatment effect to be $\tau(\mathbf{x}) := \frac{\partial}{\partial t}\mu(t, \mathbf{x})$, assuming the linearity of $\mu(t, \mathbf{x})$ with respect to t . The linearity assumption is crucial for having constant $\tau(\mathbf{x})$ for all values of t , and in this case $\tau(\mathbf{x})$ can be interpreted as the change in the average response for changing the treatment to t to $t + 1$. The linearity is a strong assumption here, and consequently, the optimal treatment will always be either the maximum or the minimum treatment value. Therefore, we avoid the use of such a technique for uplift modeling. In the following section, we propose the GCT algorithm that produces a partition Π_1, \dots, Π_K of \mathcal{X} , a partition $\Gamma_1, \dots, \Gamma_K$ of \mathcal{Z} , and the estimates of causal effects $\hat{\tau}(\Gamma_j, \Pi_i)$ for each (Γ_j, Π_i) pair. We emphasize that GCT groups similar treatments together based on homogeneity of the treatment effects jointly with the identification of homogeneous \mathbf{X} -cohorts from data. This is different from applying a multiple treatment uplift modeling after discretizing continuous treatment values based on quantiles or ad hoc cut points.

Assumptions and Limitations: Under Rubin’s potential outcome framework (Rubin, 1974), we require the following standard assumptions of the individual treatment effects (ITE) literature (Bica et al., 2020b; Imai and Van Dyk, 2004; Imbens, 2000; Schulam and Saria, 2017; van der Schaar, 2020): (i) unconfoundedness and (ii) the Stable Unit Treatment Value assumption (SUTVA). The unconfoundedness assumption states that the treatment variable and the potential outcomes are conditionally independent given the pre-treatment covariates. SUTVA ensures that the potential outcome of an experiment unit is unaffected by the treatment assignment of the other units. Additionally, we assume completely randomized treatment (i.e., the independence of the treatment variable and the pre-treatment covariates) and the existence of a control dataset (i.e., $P(T = 0) > 0$). These assumptions are more common in the uplift modeling literature than in the ITE estimation literature. It is straightforward to remove the completely randomized treatment assumption by applying,

for example, inverse propensity score weighting techniques. However, the proposed methodologies are only applicable in settings with a control dataset and a treatment dataset with different treatment values.

3 Generalized Causal Tree (GCT)

We consider data from a randomized experiment with a control group corresponding to $\{T = 0\}$ and a treatment group corresponding to $\{T \in \mathcal{Z}\}$, as described in Section 2. We define $W = \mathbb{1}_{\{T \in \mathcal{Z}\}}$ to be the treatment group indicator and $Z_1 = (T \mid W = 1)$ to be the treatment value in the treatment group. We use the subscript “1” in Z_1 to emphasize the fact that Z_1 is defined only for the treatment group $\{W = 1\}$. Let $\mathbb{S}_0^{tr} = \{(\mathbf{x}_i^{(0)}, y_i^{(0)})\}_{i=1}^{N_0}$ and $\mathbb{S}_1^{tr} = \{(\mathbf{x}_i^{(1)}, y_i^{(1)}, z_i^{(1)})\}_{i=1}^{N_1}$ denote the data corresponding to the control group $\{W = 0\}$ and the treatment group $\{W = 1\}$ respectively containing pre-treatment covariates $\mathbf{x}_i^{(w)}$ and response $y_i^{(w)}$ for $w = 0, 1$, and treatment value $z_i^{(1)}$. Note that we can fit a single treatment causal tree to \mathbb{S}^{tr} based on the binary treatment indicator W (by ignoring the Z_1 values). We generalize this by adding the use of Z_1 in the causal tree splitting criterion defined by (3.1).

Although we focus on the causal tree algorithm, our strategy works for any decision tree-based uplift modeling method. To exhibit this, we use a generic form of the objective function of the causal tree algorithm. At every step, the single treatment causal tree algorithm based on the treatment indicator W splits a leaf node in the current decision tree \mathcal{D} to increase (if possible) the value of an objective function of the form

$$\mathcal{O}(\mathcal{D}, \mathbb{S}^{tr}) = \sum_{i=1}^k g(h(\mathbb{S}_1^{tr}(\Pi_i)), h(\mathbb{S}_0^{tr}(\Pi_i))), \quad (3.1)$$

where $g(\cdot, \cdot)$ measures the utility of a leaf node (i.e., the square of estimated treatment effect minus the square of its estimated standard error), $h(\cdot)$ is a set of summary statistics (i.e., count, mean and variance) on responses conditioned on $\mathbb{S}_w^{tr}(\Pi_i) \subseteq \mathbb{S}_w^{tr}$, $\{\Pi_1, \dots, \Pi_k\}$ is a partition of \mathcal{X} defined by the leaf nodes in \mathcal{D} , and $\mathbb{S}_w^{tr}(\Pi_i) = \{(\mathbf{x}_i^{(w)}, y_i^{(w)}) : 1 \leq i \leq$

$N_w, \mathbf{x}_i^{(w)} \in \Pi_i\}$ for $w = 0, 1$. Now, we modify the objective function (3.1) to allow splitting the treatment data based on both \mathbf{X} and Z .

$$\mathcal{O}^*(\mathcal{D}, \mathbb{S}^{tr}) = \sum_{i=1}^k g(h(\mathbb{S}_1^{tr}(\Lambda_i)), h(\mathbb{S}_0^{tr}(\Pi_i))), \quad (3.2)$$

where $\mathbb{S}_1^{tr}(\Lambda_i) = \{(\mathbf{x}_i^{(1)}, y_i^{(1)}, z_i^{(1)}) : 1 \leq i \leq N_1, (\mathbf{x}_i^{(1)}, z_i^{(1)}) \in \Lambda_i\}$. We call this causal tree algorithm with the modified objective function (3.2) the basic *Generalized Causal Tree* (GCT) algorithm. The output of basic GCT is a partition $\mathcal{P}_{\mathcal{X}, \mathcal{Z}} = \{\Lambda_i\}_{i=1}^K$ of $\mathcal{X} \times \mathcal{Z}$, and the corresponding estimates $\hat{\tau}(\Lambda_i)$. Following the ‘‘honest’’ approach Athey and Imbens (2016), we assume the existence of an additional dataset $\mathbb{S}^{est} := \{\mathbb{S}_0^{est}, \mathbb{S}_1^{est}\}$ for estimating treatment effects in leaf nodes. Hence, we have

$$\hat{\tau}(\Lambda_i) = \hat{E}[Y \mid \mathbb{S}_1^{est}(\Lambda_i)] - \hat{E}[Y \mid \mathbb{S}_0^{est}(\Pi_i)],$$

i.e., the difference between the empirical average responses in $\mathbb{S}_1^{est}(\Lambda_i)$ and $\mathbb{S}_0^{est}(\Pi_i)$. In the special case where Z_1 is not chosen by GCT for partitioning (because Z_1 is not a detectable heterogeneity factor), we have $\Lambda_i = \Pi_i \times \mathcal{Z}$ and $\hat{\tau}(\Lambda_i) = \hat{\tau}(\Pi_i) = \hat{E}[Y \mid \mathbb{S}_1^{est}(\Pi_i)] - \hat{E}[Y \mid \mathbb{S}_0^{est}(\Pi_i)]$.

There are two significant obstacles to using the basic GCT algorithm in practice. The first issue is the moderately high implementation cost, as one needs to modify the splitting criterion in an existing implementation of the single treatment causal tree algorithm. The second issue is that the \mathbf{X} -cohorts $\{\Pi_i \subseteq \mathcal{X} : (\Pi_i, \Gamma_i) \in \mathcal{P}_{\mathcal{X}, \mathcal{Z}}, \Gamma_i \subseteq \mathcal{Z}\}$ are not disjoint sets and hence they do not form a partition of \mathcal{X} . This makes it computationally expensive to identify an optimal treatment for each member in a population, especially in large-scale web applications with billions of members (Tu et al., 2021). To overcome these obstacles, we propose an implementation-friendly modification of basic GCT. Then we address the second issue by making a user-friendly transformation of the output of basic GCT.

3.1 Implementation-friendly Objective

We propose a modification of the objective function (3.2) such that an existing implementation of the single treatment causal tree algorithm can be reused directly. We wish to learn a causal tree based on W while using \mathbf{X} and Z_1 jointly as heterogeneity factors. This is not possible without extending the definition of Z_1 to the control group. To this end, we assign a Z -value to each data point in the control group by randomly selecting a value from the distribution $p_{Z_1}(\cdot)$ of Z_1 . We define Z as

$$Z = Z_1 \times \mathbb{1}_{\{W=1\}} + Z_0 \times \mathbb{1}_{\{W=0\}} \quad (3.3)$$

where Z_0 is independently and identically distributed as Z_1 . In practice, we sample from the empirical distribution of Z_1 if the actual distribution is unknown. By choosing the distribution of Z_0 the same as the distribution of Z_1 , we make sure the same proportion of samples (on average) in the treatment and control group in each leaf node. Based on the definition of Z as in (3.3), we rewrite the training data as $\mathbb{S}^{tr} = \{(\mathbf{x}_i, y_i, z_i, w_i)\}_{i=1}^N$ where $N = N_0 + N_1$. Now we propose the implementation-friendly objective function

$$\mathcal{O}^*(\mathcal{D}, \mathbb{S}^{tr}) = \sum_{i=1}^k g(h(\mathbb{S}_1^{tr}(\Lambda_i)), h(\mathbb{S}_0^{tr}(\Lambda_i))), \quad (3.4)$$

where $\mathbb{S}_w^{tr}(\Lambda_i) = \{(\mathbf{x}_i, y_i, z_i, w_i) : 1 \leq i \leq N, w_i = w, (\mathbf{x}_i, z_i) \in \Lambda_i\}$ for $w = 0, 1$. Note that the difference between GCT based on Eq.(3.4) and GCT based on Eq.(3.2) is that the former splits both treatment data and control data based on a constraint on Z defined in Eq. (3.3) while the latter keeps the control data unchanged for splits based on $Z_1 = (T | W = 1)$.

Note that we reduce the effective sample size corresponding to $\{W = 0\}$ through the artificially induced randomness of Z_0 in Eq. (3.3). This could be a concern in situations where collecting data corresponding to $\{W = 0\}$ is expensive. However, this is not a concern in many cases, including web-scale problems where there is no shortage of data corresponding

to $\{W = 0\}$ (the baseline model). Typically, we have a control group at least as large as the treatment group. In the case of equal-sized treatment and control groups, GCT based on Eq. (3.4) would have identical sample size distributions for the treatment group and control group in each $\mathcal{P}_{\mathcal{X}, \mathcal{Z}}$ cohort. This implies the same level of uncertainties in the average treatment response and average control response in those groups.

Additional Notation: Consider a decision tree \mathcal{D} based on a features \mathbf{X} . Recall that \mathcal{D} defines a partition of \mathcal{X} based on its leaf nodes. We assign a unique identifier $id(\ell, \mathcal{D})$ to each leaf ℓ . We represent \mathcal{D} as a full binary tree where each non-leaf node a contains (i) a pointer to the left child node $left(a)$, (ii) a pointer to the right child node $right(a)$, and (iii) a cohort definition $\pi(a) \subseteq \mathcal{X}$ with $\pi(a) := \mathcal{X}$ for the root node a . A non-leaf node a is a unique parent of $left(a)$ and $right(a)$. If we decide to split a node a based on a continuous feature $X_r \in \mathbf{X}$ and at the value s , then we have $\pi(left(a)) = \{(x_r, \mathbf{x}_{-r}) \in \mathcal{X}_r \times \mathcal{X}_{-r} : x_r \leq s\}$ and $\pi(right(a)) = \{(x_r, \mathbf{x}_{-r}) \in \mathcal{X}_r \times \mathcal{X}_{-r} : x_r > s\}$, where \mathcal{X}_r and \mathcal{X}_{-r} denote the sets of values taken by X_r and $\mathbf{X} \setminus \{X_r\}$ respectively. We denote the parent of a in \mathcal{D} by $parent(a, \mathcal{D})$. A sequence of nodes a_0, \dots, a_k is an ancestral path a_0 to a_k in \mathcal{D} if for each $i \in \{1, \dots, k\}$, $a_{i-1} = parent(a_i, \mathcal{D})$. Let $\{\ell_i\}_{i=1}^K$ be all leaf nodes of \mathcal{D} . For each leaf ℓ_i , we define $\Pi_i := \Pi_{\mathbf{X}}(\ell_i, \mathcal{D}) := \bigcap_{r=0}^{k_i} \pi(a_r, \mathcal{D})$, where a_0, \dots, a_{k_i} is the unique ancestral path from the root a_0 to $\ell_i = a_{k_i}$. Then $\mathcal{P}_{\mathcal{X}} := \{\Pi_i\}_{i=1}^K$ forms a partition of \mathcal{X} .

A post-order traversal on a binary tree returns an ordering of the node through a recursive depth-first traversal that first visits the left sub-tree, then the right sub-tree, and finally the root. For any $\mathbf{X}' \subseteq \mathbf{X}$, we denote the coordinate-wise projection of $\Pi_{\mathbf{X}}(a, \mathcal{D})$ onto \mathcal{X}' (the set of all possible \mathbf{X}' values) by $\Pi_{\mathbf{X}'}(a, \mathcal{D})$ (i.e., obtained by removing all coordinates corresponding to $\mathbf{X} \setminus \mathbf{X}'$). For convenience, we use the same notation to denote a singleton set $\{i\}$ and its element i .

3.2 User-friendly Output

We transform the output of basic GCT $\mathcal{P}_{\mathcal{X}, \mathcal{Z}}$ to a cross-product $\mathcal{P}_{\mathcal{X}} \times \mathcal{P}_{\mathcal{Z}} = \{\Pi_1, \dots, \Pi_{K_{\mathcal{X}}}\} \times \{\Gamma_1, \dots, \Gamma_{K_{\mathcal{Z}}}\}$ of partitions of \mathcal{X} and partitions of \mathcal{Z} and provide corresponding treatment

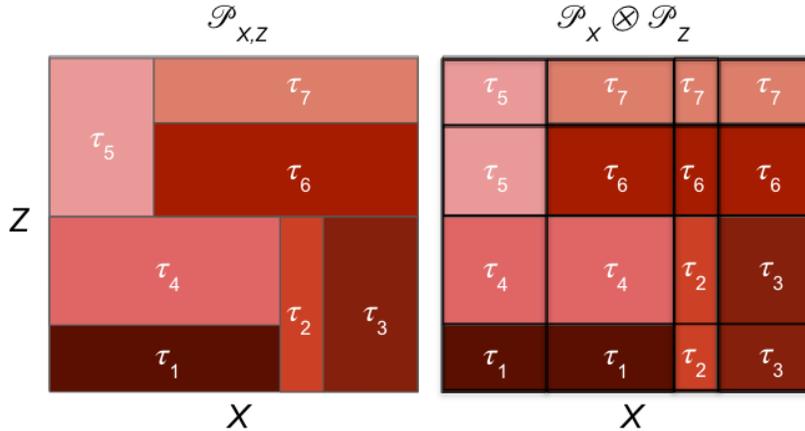


Figure 1: GCT output transformation visualization

effect estimates, i.e., $\{\hat{\tau}(\Pi_i, \Gamma_j) : \Pi_i \in \mathcal{P}_X, \Gamma_j \in \mathcal{P}_Z\}$. The transformed output is much more user-friendly because the cross-product representation of cohorts makes it easy to find an optimal treatment set for each $\Pi_i \in \mathcal{P}_X$ by setting $\Gamma_{i*} = \operatorname{argmax}_{\Gamma_j \in \mathcal{P}_Z} \hat{\tau}(\Pi_i, \Gamma_j)$ whenever $\max_{\Gamma_j \in \mathcal{P}_Z} \hat{\tau}(\Pi_i, \Gamma_j) > 0$. Then for any individual with features $\mathbf{x} \in \Pi_i$ an optimal treatment can be chosen from the conditional distribution $p_Z(\cdot | \Gamma_{i*})$ of Z given $\{Z \in \Gamma_{i*}\}$.

We use tree-modification procedure (Algorithm 1) to obtain $\mathcal{P}_X \times \mathcal{P}_Z$ and the corresponding treatment effects via Algorithm 5. Before that, we illustrate the main idea with the toy example depicted in Figure 1. We consider a continuous-valued Z and an one-dimensional continuous-valued X . The left sub-figure represents estimated causal effects in a basic GCT as piecewise constant functions on $\mathcal{P}_{X,Z}$, which is a collection of disjoint and exhaustive axis-aligned rectangles. The right sub-figure represents $\mathcal{P}_X \times \mathcal{P}_Z$ as a minimal (in terms of the total number of rectangles) finer partition of $\mathcal{X} \times \mathcal{Z}$ that is a cross-product of a partition of \mathcal{X} (corresponding to the vertical lines in the sub-figure) and a partition of \mathcal{Z} (corresponding to the horizontal lines in the sub-figure). The (color-coded) treatment effects in these finer partitions are borrowed from the treatment effects in $\mathcal{P}_{X,Z}$.

Note that the X -cohorts can be obtained by projecting each cohort $\Lambda \in \mathcal{P}_{X,Z}$ onto \mathcal{X} and by splitting the resulting non-disjoint sets into disjoint sets appropriately. The splitting steps can be done efficiently by exploiting the underlying decision tree structure. We present

this idea formally in Algorithm 1, which allows us to remove a set of features \mathbf{X}' from a decision tree \mathcal{D}_{in} such that (1) the output is a decision tree \mathcal{D}_{out} based on the feature set $\mathbf{X} \setminus \mathbf{X}'$, and (2) for each leaf node ℓ in \mathcal{D}_{out} , there exists a leaf node ℓ' in \mathcal{D}_{in} such that $\Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell, \mathcal{D}_{out}) \subseteq \Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell', \mathcal{D}_{in})$, where $\Pi_{\mathbf{X} \setminus \mathbf{X}'}(\cdot, \mathcal{D}_{in})$ denotes a cohort of \mathcal{D}_{in} projected onto $\mathcal{X} \setminus \mathcal{X}'$. Additionally, For each leaf ℓ in \mathcal{D}_{out} , we facilitate the identification of the set of super-cohort leaves in \mathcal{D}_{in} , defined as

$$\{\ell' \in \mathcal{D}_{in} : \Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell, \mathcal{D}_{out}) \subseteq \Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell', \mathcal{D}_{in})\}. \quad (3.5)$$

Algorithm 1 Removing features from a decision tree

Require: a decision tree \mathcal{D} , a set features \mathbf{X}'

- 1: Let a_1, \dots, a_n be the ordering of the nodes from a post-order traversal in \mathcal{D} ;
 - 2: **for** $i = 1, \dots, n$ **do**
 - 3: **if** a_i is a non-leaf and $\pi(left(a))$ (or $\pi(right(a))$) contains a constraint involving a feature in \mathbf{X}' , **then**
 - 4: merge $left(a_i)$ and $right(a_i)$ as follows:
 - 5: **if** both $left(a_i), right(a_i)$ are leaves **then**
 - 6: MergeLeaves(\mathcal{D}, a_i);
 - 7: **else if** one of $left(a_i), right(a_i)$ is a leaf **then**
 - 8: MergeLeafAndNonLeaf(\mathcal{D}, a_i);
 - 9: **else**
 - 10: MergeNonLeaves(\mathcal{D}, a_i);
 - 11: **end if**
 - 12: **end if**
 - 13: **end for**
 - 14: Return \mathcal{D}
-

These sets of super-cohort leaves will be used in Algorithm 5 for computing treatment effects. The main idea of Algorithm 1 is to sequentially visit each node a via a post-order traversal and remove its children ($left(a, \mathcal{D}), right(a, \mathcal{D})$) whenever the split is defined by a feature in \mathbf{X}' , and merge the sub-trees below the children with the parent node a . The merging operation involves updating the children pointers and removing all nodes corresponding to a conflicting cohort definition. Moreover, we keep track of the super-cohort sets defined in (3.5) by passing the identifier information of the deleted nodes and the nodes that has

lost their leaf status due to the merging.

The removing algorithm sequentially visits each node via a post-order traversal and removes a split from the tree whenever the split is associated with the targeted feature. Once a split is identified as a to-be-removed branch, the removal process not only concerns with taking out the node from the tree, but also requires attention on merging the sub-trees below the identified node with its parent node. Let a denote the parent node, and we want to merge the two children of node a . Depending on the positions of the two child nodes, i.e., $left(a)$ and $right(a)$, the merging operations can be categorized into three cases:

- (i) **MergeLeaves** (Algorithm 2): When both children are leaf nodes, we merge the two by removing them from the tree, and by passing their identifier to their parents, i.e., $id(a, \mathcal{D}) = id(left(a), \mathcal{D}) \cup id(right(a), \mathcal{D})$. See Figure (2a) for a graphical illustration.
- (ii) **MergeLeafAndNonLeaf** (Algorithm 3): When one child of a is leaf, denoted by c_2 , and the other child is not a leaf, denoted by c_1 , we connect the sub-tree below c_1 with the parent of c_1 . In this case, the identifier of c_2 gets passed to the leaf nodes below c_1 . See Figure (2b).
- (iii) **MergeNonLeaves** (Algorithm 4): When both children are not leaves, denoted by c_1 and c_2 , we first connect the sub-tree below c_1 with the parent of c_1 . Then we connect the sub-tree below c_2 with each leaf below c_1 (node B , C and $right(c_1)$ in Figure (2c)) while passing the identifiers to the leaves in the merged tree. This merging can cause some empty cohort due to conflicting constraints and we prune the tree by removing nodes containing empty cohorts. See Fig (2c) for an example.

Note that whenever we remove a leaf node ℓ from the current tree $\mathcal{D}_{current}$ in **MergeLeaves** and **MergeLeafAndNonLeaf** or transform a leaf node ℓ to non-leaf node in $\mathcal{D}_{current}$ in **MergeNonLeaves**, we make to pass the identifiers to all leaf nodes ℓ' in resulting tree \mathcal{D}_{new} such that $\Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell', \mathcal{D}_{new}) \subseteq \Pi_{\mathbf{X} \setminus \mathbf{X}'}(\ell, \mathcal{D}_{current})$. This facilitates the retrieval of S_ℓ defined in (3.5) from the output \mathcal{D}_{out} of Algorithm 1 as $S_\ell = \{\ell' \in \mathcal{D}_{in} : id(\ell', \mathcal{D}_{in}) \in id(\ell, \mathcal{D}_{out})\}$.

Algorithm 2 MergeLeaves(\mathcal{D}, a)

Require: a decision tree \mathcal{D} , a node a such that both children of a are leaf nodes

- 1: Set $id(a, \mathcal{D}) = id(left(a), \mathcal{D}) \cup id(right(a), \mathcal{D})$;
 - 2: Set $left(a) = right(a) = \emptyset$;
 - 3: Return \mathcal{D} ;
-

Algorithm 3 MergeLeafAndNonLeaf(\mathcal{D}, a)

Require: a decision tree \mathcal{D} , a node a such that one child of a is a leaf and the other child is not a leaf

- 1: Let c_1 be the non-leaf child of a and let c_2 be the leaf child;
 - 2: Set $left(a, \mathcal{D}) = left(c_1)$ and $right(a) = right(c_1)$;
 - 3: **for** each visited node a_1 in a breadth-first traversal in the sub-tree below a **do**
 - 4: **if** a_1 is a leaf node **then**
 - 5: Set $id(a_1, \mathcal{D}) = id(a_1, \mathcal{D}) \cup id(c_2, \mathcal{D})$;
 - 6: **end if**
 - 7: **end for**
 - 8: Return \mathcal{D} ;
-

Algorithm 4 MergeNonLeaves(\mathcal{D}, a)

Require: a decision tree \mathcal{D} , a node a such that both children of a are non-leaf nodes

- 1: Let $c_1 = left(a)$ and $c_2 = right(a)$;
 - 2: Set $left(a) = left(c_1)$ and $right(a) = right(c_1)$;
 - 3: **for** each visited node a_1 in a breadth first traversal in the sub-tree below a **do**
 - 4: **if** a_1 is a leaf node **then**
 - 5: Set $left(a_1, \mathcal{D}) = left(c_2, \mathcal{D})$ and $right(a_1, \mathcal{D}) = right(c_2, \mathcal{D})$;
 - 6: **for** each visited node a_2 in a breadth first traversal in the sub-tree starting from a_1 **do**
 - 7: **if** a_2 is a leaf node **then**
 - 8: $id(a_2, \mathcal{D}) = id(a_1, \mathcal{D}) \cup id(a_2, \mathcal{D})$;
 - 9: **else**
 - 10: **if** $\Pi_{\mathbf{X}}(c, \mathcal{D}) = \emptyset$ for a child of a_2 **then**
 - 11: Let c' be the other child;
 - 12: Remove c and c' by setting $left(a_2) = left(c')$ and $right(a_2) = right(c')$ and
 - 13: **if** c' is a leaf **then**
 - 14: $id(a_2, \mathcal{D}) = id(c', \mathcal{D})$;
 - 15: **end if**
 - 16: **end if**
 - 17: **end if**
 - 18: **end for**
 - 19: **end if**
 - 20: **end for**
 - 21: **end for**
 - 21: Return \mathcal{D} ;
-

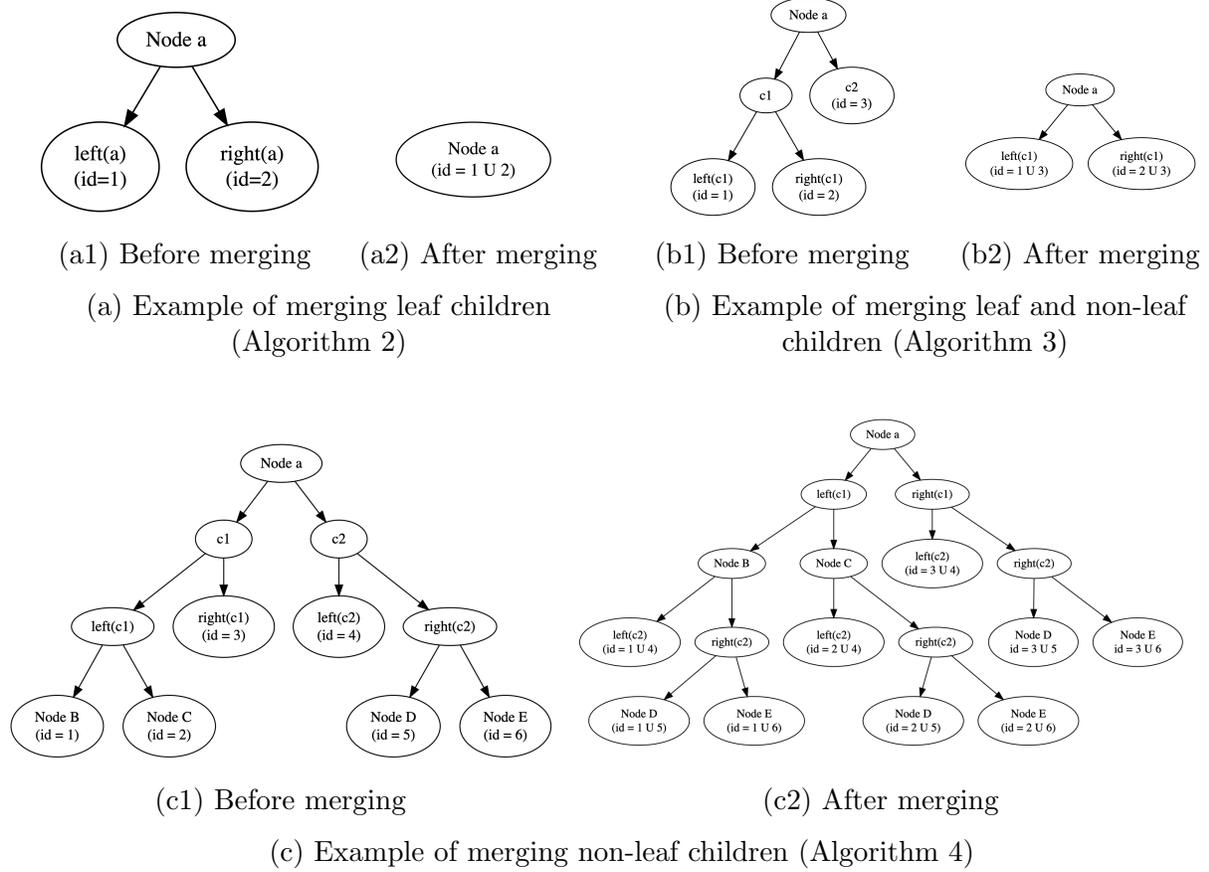


Figure 2: Examples of Algorithms 2-4

3.3 The Complete Algorithm

We formally present the GCT algorithm for identifying heterogeneous cohorts and estimating the treatment effects of those cohorts in Algorithm 5. The algorithm is based on the implementation-friendly objective function defined in Section 3.1 and user-friendly output given in Section 3.2.

The output of the algorithm is a partition $\{\Pi_1, \dots, \Pi_{K_{\mathcal{X}}}\}$ of \mathcal{X} , a partition $\{\Gamma_1, \dots, \Gamma_{K_{\mathcal{Z}}}\}$ of \mathcal{Z} , and the treatment effects $\hat{\tau}(\Pi_i, \Gamma_j)$ for $i = 1, \dots, K_{\mathcal{X}}$ and $j = 1, \dots, K_{\mathcal{Z}}$. Given these estimates, for any $x \in \mathcal{X}$ and $z \in \mathcal{Z}$, we can estimate the treatment effect as $\hat{\tau}(\Pi_i, \Gamma_j)$ by identifying the unique cohorts Π_i and Γ_j such that $x \in \Pi_i$ and $z \in \Gamma_j$. A step-by-step demonstration of Algorithm 5 is provided in Appendix A.

Algorithm 5 Generalized Causal Tree

- 1: Learn a causal tree $\mathcal{C}_{\mathbf{X},Z}$ with treatment W , heterogeneity factors (\mathbf{X}, Z) and response Y ;
 - 2: Let $\Lambda_1, \dots, \Lambda_M$ be the leaf-cohorts of $\mathcal{C}_{\mathbf{X},Z}$ corresponding to leaf nodes ℓ_1, \dots, ℓ_M ;
 - 3: Let $\hat{\tau}(\Lambda_i) := \hat{E}[Y(W = 1)|(\mathbf{X}, Z) \in \Lambda_i] - \hat{E}[Y(W = 0)|(\mathbf{X}, Z) \in \Lambda_i]$ for $i = 1, \dots, M$;
 - 4: Remove all nodes involving Z from $\mathcal{C}_{\mathbf{X},Z}$ using Algorithm 1 with $\mathcal{D} = \mathcal{C}_{\mathbf{X},Z}$ and $\mathbf{X}' = \{Z\}$ and denote the output by $\mathcal{C}_{\mathbf{X}}$ with leaf-cohorts $\Pi_1, \dots, \Pi_{K_{\mathbf{X}}}$ corresponding to leaf nodes $\ell_{\mathbf{X},1}, \dots, \ell_{\mathbf{X},K_{\mathbf{X}}}$;
 - 5: Remove all nodes involving \mathbf{X} from $\mathcal{C}_{\mathbf{X},Z}$ using Algorithm 1 with $\mathcal{D} = \mathcal{C}_{\mathbf{X},Z}$ and $\mathbf{X}' = \mathbf{X}$ and denote the output by \mathcal{C}_Z with leaf-cohorts $\{\Gamma_1, \dots, \Gamma_{K_Z}\}$ corresponding to leaf nodes $\ell_{Z,1}, \dots, \ell_{Z,K_Z}$;
 - 6: **for** each (Π_i, Γ_j) **do**
 - 7: Find the unique leaf node ℓ in $\mathcal{C}_{\mathbf{X},Z}$ corresponding to the id $id(\ell_{\mathbf{X},i}, \mathcal{C}_{\mathbf{X}}) \cap id(\ell_{Z,j}, \mathcal{C}_Z)$;
 - 8: Set $\hat{\tau}(\Pi_i, \Gamma_j) = \hat{\tau}(\Lambda_{\ell})$;
 - 9: **end for**
-

Theoretical Properties: Our algorithm is designed in a way such that all statistical properties of the single treatment causal tree based treatment effect estimators hold for our GCT based treatment effect estimators. This follows from the following identifiability result where we show that the causal tree algorithm based on the treatment indicator W and the heterogeneity factors (X, \mathbf{Z}) is indeed estimating $\tau(z, \mathbf{x})$. The proof is given in Appendix B.

Theorem 1. *Let $\tau(z, \mathbf{x})$ be the conditional average treatment effect of $T = z$ given $\mathbf{X} = x$. Then $\tau(z, \mathbf{x}) = E[Y(W = 1) | \mathbf{X} = \mathbf{x}, Z = z] - E[Y(W = 0) | \mathbf{X} = \mathbf{x}, Z = z]$ under the assumptions given in Section 1.*

Scalability: Since we are just adding one extra variable Z to the feature set \mathbf{X} , the impact on the time-complexity of the original causal-tree algorithm (or any decision tree based method) is expected to be minimal. It is easy to show that the worst-case time-complexity of creating the used-friendly output from a given causal tree with m nodes is $\mathcal{O}(m^2)$.

4 Experiments

We conduct simulation studies to demonstrate the efficacy of GCT for handling multiple discrete or continuous treatments. Here, we consider (i) continuous treatment: Z following a

uniform distribution on $\mathcal{Z} = (0, 1]$, (ii) ordinal treatment: Z following a uniform distribution on the ordered set $\mathcal{Z} = \{1, 2, 3, 4, 5, 6\}$, and (iii) categorical treatment: Z follows a uniform distribution on the set $\mathcal{Z} = \{a, b, c, d\}$. Note that the GCT algorithm treats in the ordinal and categorical treatment cases slightly differently. In the ordinal case, GCT anticipates the function $\tau(\mathbf{x}, t)$ to be a smooth function of t for each \mathbf{x} and considers splitting the treatment values $\{1, 2, \dots, M\}$ into two groups (while building the tree) based on $M - 1$ choices (i.e., $\{1, \dots, i\}$ and $\{i + 1, \dots, M\}$ for $i = 1, \dots, M - 1$) instead of all 2^M choices (e.g., the split $\{1, M\}$ and $\{2, \dots, M - 1\}$ is not considered in the ordinal case).

For each setup, we generate a sample of $N = 2000$ individuals with heterogeneity features $\mathbf{X} = (X_1, X_2)$, where X_1 and X_2 are independently generated from a *Uniform*(0, 1) distribution. The sample is randomly splitted into two halves: one for training (denoted by $\{\mathbf{X}^{tr}\}$) and one for test (denoted by $\{\mathbf{X}^{te}\}$). On the training data, we independently generate the treatments T_i^{tr} for each individual i . Conditioning on $(\mathbf{X}_i^{tr}, T_i^{tr})$, the outcome Y_i^{tr} is generated according to $Y_i^{tr} = f(\mathbf{X}_i^{tr}, T_i^{tr}) + \epsilon_i^{tr}$, where ϵ_i^{tr} 's are i.i.d. $\mathcal{N}(0, 1)$ random variables and $f(\cdot, \cdot)$ is defined below. From the training data, we learn \mathbf{X} -cohorts (Π_1, \dots, Π_K) along with the within-cohort optimal treatment set Γ_{j_r} for each Π_r . Then, on the test data, for each \mathbf{X}_i^{te} , we determine the cohort Π_r such that $\mathbf{X}_i^{te} \in \Pi_r$ and generate T_i^{te} by randomly choosing a value from the uniform distribution on Γ_{j_r} . Conditioning on $(\mathbf{X}_i^{te}, T_i^{te})$, the outcome Y_i^{te} is generated according to $Y_i^{te} = f(\mathbf{X}_i^{te}, T_i^{te}) + \epsilon_i^{te}$, where ϵ_i^{te} 's are i.i.d. $\mathcal{N}(0, 1)$ random variables. We define $f(\cdot, \cdot)$ through the function $\eta(x_1, x_2) = -2 + 4/([1 + \exp(-12(x_1 - 0.2))][1 + \exp(-12(x_2 - 0.2))])$. of heterogeneous features $\mathbf{X}_i = (X_{1i}, X_{2i})$. Figure 3 in the appendix provides a graphical illustration of the values of $\eta(x_1, x_2)$ for (x_1, x_2) within a unit square. Blue and red depict the regions in which $\eta(x_1, x_2)$ takes positive and negative values, respectively. The heterogeneous nature of how the function $\eta(x_1, x_2)$ responds to features makes it the building block for designing our experiments. Besides $\eta(x_1, x_2)$, we include $\eta(x_1, 1 - x_2)$ to infuse more heterogeneity to the data. We define different $f(\cdot, \cdot)$ functions corresponding to the different types of treatments

such as continuous, ordinal, categorical. In Appendix C, we provide discussions on the choice and the intuition behind our simulation designs.

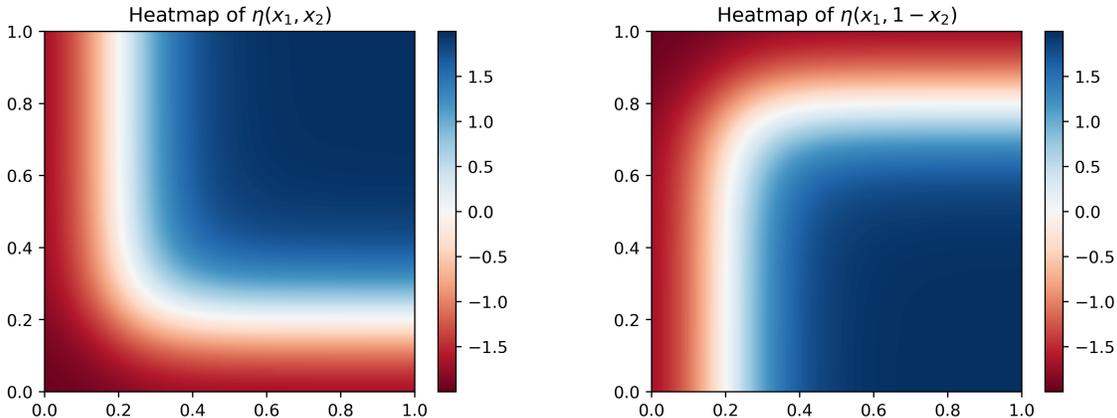


Figure 3: Heatmaps of $\eta(x_1, x_2)$ and $\eta(x_1, 1 - x_2)$.

We compare GCT¹ with existing decision tree based uplift modeling, including the original causal tree Athey and Imbens (2016) and multi-treatment uplift models of Radcliffe and Surry (2011); Zhao et al. (2017) based on various split criteria including the KL divergence (KL), Euclidean distance (ED), Chi-Square (CHI), and contextual treatment selection (CTS)². Since existing methods are not defined for continuous treatments, we discretize the continuous treatment into 4 equally sized bins for these methods and use the mid-point of each bin to represent the treatment level for that bin. Recall that the original causal tree algorithm can only work with a single treatment. Therefore, when learning the heterogeneous cohorts using CT, we have to ignore the treatment values Z but only use the treatment indicator W . The resultant estimated causal effects for each cohort are computed on a treatment-vs-control basis regardless of treatment values. To assign the optimal treatment, we randomly choose a value from a uniform distribution on \mathcal{Z} whenever the treatment effect is positive. We refer to this variation of causal tree by CT-B (where B stands for “binary”). We consider another variation, called CT-M (M for “multiple”), where we estimate causal

¹The code is available through [Github](#).

²CT is implemented using “causalTree” R package, and the latter four methods are implemented using “CausalML” Python package (Chen et al., 2020).

effects for every discrete (or discretized) treatment level separately within each cohort after applying the causal tree based on W , and we assign optimal treatments based on these estimates.

The algorithm performance is evaluated by two metrics: (1) the mean square error (MSE) of estimated causal effects, and (2) the average outcome with the optimal treatment allocated to the testing data. The former measures the accuracy of an algorithm in estimating heterogeneous treatment effects, while the latter measures uplifting quality based on the optimal treatment allocation according to different algorithms. Figure 4 shows box-plots of MSE values. Results are summarized in Table 1. Each metric is averaged over 100 replications. The standard errors are reported in the parentheses.

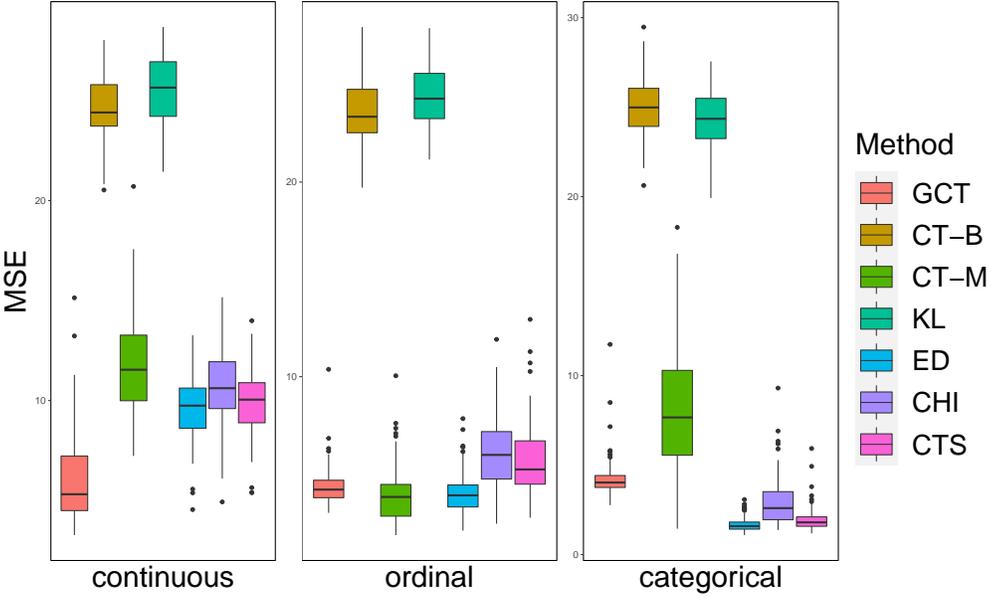


Figure 4: Box-plots of the MSE of estimated treatment effects based on 100 iterations.

From Figure 4, we can see that GCT achieves the smallest MSE value in the continuous setting and yields comparable MSE with ED, CHI, and CTS in the ordinal and categorical settings. More importantly, Table 1 shows that we can achieve the highest average outcome by determining the optimal treatment allocations based on the output of GCT. The negative outcome obtained by CT-B in the categorical setting further implies that neglecting the

existence of heterogeneous effects among various treatments may lead to unfavorable results.

Table 1: Average outcome with optimal treatment allocations (a higher value implies a better performance).

Method	Treatment Type		
	Continuous	Ordinal	Categorical
GCT	5.744 (0.051)	5.321 (0.048)	4.480 (0.045)
CT-B	0.188 (0.013)	0.581 (0.017)	-0.008 (0.012)
CT-M	4.880 (0.092)	3.904 (0.033)	3.618 (0.057)
KL	2.519 (0.033)	2.519 (0.033)	2.502 (0.039)
ED	0.548 (0.036)	0.417 (0.041)	0.466 (0.033)
CHI	0.698 (0.038)	0.843 (0.050)	0.626 (0.037)
CTS	0.598 (0.035)	0.818 (0.050)	0.510 (0.035)

Table 2: Expected response $\bar{R} = \frac{1}{N} \sum_{i=1}^N R_i$.

Dataset	GCT	CT-B	CT-M	KL	ED	CHI	CTS
Bladder	0.81	0.76	0.79	0.82	0.81	0.82	0.81
AOD	0.62	0.58	0.60	0.61	0.60	0.56	0.59

5 Real Data Analysis

In a recent benchmarking study of multi-treatment uplift modeling methods (Olaya et al., 2020), the authors compared tree-based methods with other (less-interpretable) methods to conclude that no approach is strictly superior to the others. We used two datasets from the same benchmarking setup to compare our method with the tree-based methods and we refer to Olaya et al. (2020) for a comparison with the other multi-treatment uplift modeling methods. In particular, we apply the proposed GCT as well as CT-B, CT-M, and the four tree-based methods KL, ED, CHI, and CTS to two real datasets to evaluate their uplifting effects. One is the *Bladder* dataset on the recurrence of bladder cancer with weak heterogeneous effect, available in the R package `survival` (Therneau, 2021). The other is the *AOD* dataset on alcohol and drug usage with a strong heterogeneous effect, available in the R package `twang` (Cefalu et al., 2021). There are 2 treatment groups and 1 control group in both datasets and the outcome for Bladder is binary while the outcome for AOD

is continuous. Following Olaya et al. (2020), we convert the continuous outcome to a binary outcome in AOD for a fair comparison. Unlike the simulation studies, it is impossible to re-run the experiments based on the optimal treatment determined by our algorithm. Instead, we follow the proposal of Zhao et al. (2017) for quantifying the performance of an algorithm. The individual expected response is defined as $R_i = \sum_{k=1}^K \frac{y_i}{P_k} \mathbb{1}_{\{v(\mathbf{x}_i)=k\}} \mathbb{1}_{\{T_i=k\}}$, where P_k is the empirical probability of an individual receiving treatment k , and $v(\mathbf{x}_i)$ is the optimal treatment assigned by the uplift model based on its heterogeneous features \mathbf{x}_i . The expected response of an uplift model is then defined as $\bar{R} = \frac{1}{N} \sum_{i=1}^N R_i$, which is an unbiased estimator of $E[y_i|T = v(\mathbf{x}_i)]$. Following Olaya et al. (2020), we first apply propensity scoring matching to debias the non-experimental data. Then we fit models using 5-fold cross-validation, and the results from each round are averaged to obtain the overall performance.

Table 2 summarizes the average expected responses for the 7 methods. Notice that results of KL and ED are not exactly the same as but very close to those in Olaya et al. (2020) due to the randomness in the training-test split. For both datasets, all methods performed similarly. This was expected as the number of treatments is small (i.e., two). GCT is expected to be more beneficial with a higher number of treatments (or continuous treatments), as we have shown in simulations.

6 Discussion

We presented a novel generalization of the causal tree algorithm for uplift modeling with multiple discrete or continuous treatments. One of the key features of our generalized causal tree (GCT) algorithm is that it provides a data-driven way of grouping similar treatments together in addition to partitioning the underlying population into disjoint cohorts based on the homogeneity of the treatment effects. We provide significant improvements over the basic version of GCT to make it implementation-friendly and make its output easily adoptable to downstream applications, such as optimal treatment allocation. Note that the implementation-friendly modification might lead to performance loss by creating unnecessary

splits in the control data, which is expected to be negligible when the size of the control data is much larger than the training data.

We empirically showed that GCT outperforms some extensions of the original causal tree algorithm as well as existing decision tree based multiple treatment uplift modeling methods. Note that our generalization strategy is not limited to the splitting criterion of the causal tree algorithm and hence can be used to extend any decision tree based uplift modeling technique to multiple discrete or continuous treatments. We acknowledge that the decision tree based approaches may not be the best performing methods in certain contexts. However, they might still be preferable because of their interpretability and the computational efficiency in using their output in optimal treatment allocation. Finally, we also note that GCT can be directly used in a continuous treatment version of some recent extensions to uplift modeling, such as treatments with different costs (Zhao and Harinen, 2019), and constrained utility optimization (Tu et al., 2021).

References

- Agarwal, D., Basu, K., Ghosh, S., Xuan, Y., Yang, Y., and Zhang, L. (2018). Online parameter selection for web-based ranking problems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 23–32.
- Athey, S. and Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360.
- Bica, I., Jordon, J., and van der Schaar, M. (2020a). Estimating the effects of continuous-valued interventions using generative adversarial networks. *Advances in Neural Information Processing Systems*, 33:16434–16445.
- Bica, I., Jordon, J., and van der Schaar, M. (2020b). Estimating the effects of continuous-valued interventions using generative adversarial networks. In *Proceedings of the 34th*

- International Conference on Neural Information Processing Systems*. Curran Associates Inc.
- Cefalu, M., Ridgeway, G., McCaffrey, D., Morral, A., Griffin, B. A., and Burgette, L. (2021). *twang: Toolkit for Weighting and Analysis of Nonequivalent Groups*. R package version 2.4.
- Chen, H., Harinen, T., Lee, J.-Y., Yung, M., and Zhao, Z. (2020). CausalML: Python package for causal machine learning. *arXiv preprint arXiv:2002.11631*.
- Curth, A. and van der Schaar, M. (2021). Nonparametric estimation of heterogeneous treatment effects: From theory to learning algorithms. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics*.
- Green, D. P. and Kern, H. L. (2012). Modeling Heterogeneous Treatment Effects in Survey Experiments with Bayesian Additive Regression Trees. *Public Opinion Quarterly*, 76(3):491–511.
- Guelman, L., Guillén, M., and Pérez-Marín, A. M. (2014). Optimal personalized treatment rules for marketing interventions: A review of methods, a new proposal, and an insurance case study.
- Guo, R., Cheng, L., Li, J., Hahn, P. R., and Liu, H. (2020). A survey of learning causality with data: Problems and methods. *ACM Comput.Surv.*, 53(4).
- Gupta, R., Liang, G., Tseng, H.-P., Holur Vijay, R. K., Chen, X., and Rosales, R. (2016). Email volume optimization at linkedin. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 97–106.
- Gutierrez, P. and Gérardy, J.-Y. (2017). Causal inference and uplift modelling: A review of the literature. In *International Conference on Predictive Applications and APIs*, pages 1–13. PMLR.

- Hansotia, B. and Rukstales, B. (2002). Incremental value modeling. *Journal of Interactive Marketing*, 16:35–46.
- Hill, J. L. (2011). Bayesian nonparametric modeling for causal inference. *J Comput Graphical Stat*, 20:217–240.
- Imai, K. and Van Dyk, D. A. (2004). Causal inference with general treatment regimes: Generalizing the propensity score. *Journal of the American Statistical Association*, 99(467):854–866.
- Imbens, G. (2000). The role of the propensity score in estimating dose-response functions. *Biometrika*, 87(3):706–710.
- Jin, H. and Rubin, D. B. (2008). Principal stratification for causal inference with extended partial compliance. *Journal of the American Statistical Association*, 103(481):101–111.
- Knaus, M. C., Lechner, M., and Strittmatter, A. (2020). 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., and Yu, B. (2019). Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116:4156–4165.
- Liu, W., Yu, X., Mao, J., Wu, X., and Dyer, J. (2023). Quantifying the effectiveness of advertising: A bootstrap proportion test for brand lift testing. In *CIKM’23*.
- Olaya, D., Coussement, K., and Verbeke, W. (2020). A survey and benchmarking study of multitreatment uplift modeling. *Data Mining and Knowledge Discovery*, 34(2):273–308.
- Oprescu, M., Syrgkanis, V., and Wu, Z. S. (2019). Orthogonal random forest for causal inference. In *International Conference on Machine Learning*, pages 4932–4941. PMLR.

- Radcliffe, N. and Surry, P. (2011). Real-world uplift modelling with significance-based uplift trees. *Stochastic Solutions White Paper*.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688.
- Rzepakowski, P. and Jaroszewicz, S. (2012). Decision trees for uplift modeling with single and multiple treatments. *Knowledge and Information Systems*, 32(2):303–327.
- Schulam, P. and Saria, S. (2017). Reliable decision support using counterfactual models. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, page 1696–1706, Red Hook, NY, USA. Curran Associates Inc.
- Therneau, T. M. (2021). *A Package for Survival Analysis in R*. R package version 3.2-13.
- Tu, Y., Basu, K., DiCiccio, C., Bansal, R., Nandy, P., Jaikumar, P., and Chatterjee, S. (2021). Personalized treatment selection using causal heterogeneity. In *Proceedings of The Web Conference 2021*, WWW ’21.
- van der Schaar, Y. Z. A. B. M. (2020). Learning overlapping representations for the estimation of individualized treatment effects. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242.
- Wan, S., Zheng, C., Sun, Z., Xu, M., Yang, X., Zhu, H., and Guo, J. (2022). Gcf: Generalized causal forest for heterogeneous treatment effect estimation in online marketplace. *arXiv preprint arXiv:2203.10975*.
- Zaniewicz, L. and Jaroszewicz, S. (2013). Support vector machines for uplift modeling. In *13th International Conference on Data Mining Workshops*, pages 131–138.

- Zhang, W., Li, J., and Liu, L. (2021). A unified survey of treatment effect heterogeneity modelling and uplift modelling. *ACM Computing Surveys (CSUR)*, 54(8):1–36.
- Zhao, Y., Fang, X., and Simchi-Levi, D. (2017). Uplift modeling with multiple treatments and general response types. In *Proceedings of the 2017 SIAM International Conference on Data Mining*, pages 588–596. SIAM.
- Zhao, Z. and Harinen, T. (2019). Uplift modeling for multiple treatments with cost optimization. In *2019 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 422–431. IEEE.
- Zhou, Z., Athey, S., and Wager, S. (2023). Offline multi-action policy learning: Generalization and optimization. *Operations Research*, 71(1):148–183.

Appendices

The appendices provide supplemental details on the algorithms and numerical results to Sections 3-4 of the main context. Section A uses a toy example to provide a step-by-step demonstration of the implementation of our proposed GCT algorithm. Section B presents proofs. Section C discusses the simulation designs.

A A Toy Example

This section presents a toy example of our proposed GCT approach. Figure 5 plots a causal tree $\mathcal{C}_{\mathbf{X},Z}$ with respect to heterogeneity factors $\mathbf{X} = (X_1, X_2)$ and treatment values Z .

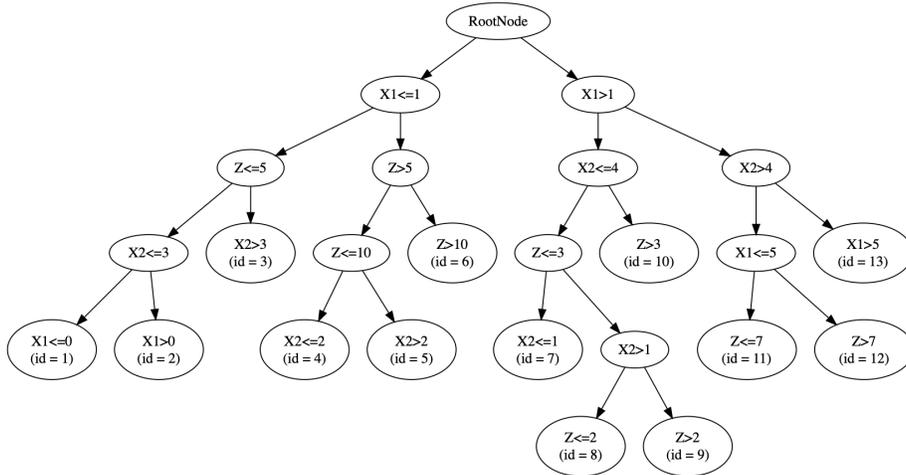


Figure 5: An example of causal tree $\mathcal{C}_{\mathbf{X},Z}$

From $\mathcal{C}_{\mathbf{X},Z}$, we obtain the leaf-cohorts $\Lambda_1, \dots, \Lambda_M$ with $M = 13$. Each leaf-cohort corresponds to a leaf node in the tree. To be more specific, we summarize the details of all leaf cohorts in Table 3.

Table 3: Leaf-cohorts $\Lambda_1, \dots, \Lambda_M$ from the causal tree $\mathcal{C}_{\mathbf{X},Z}$

Leaf-cohorts	Leaf nodes	Cohort regions	Leaf-cohorts	Leaf id	Cohort regions
Λ_1	ℓ_1 : id = 1	$\{X_1 \leq 0, X_2 \leq 3, Z \leq 5\}$	Λ_8	ℓ_8 : id = 8	$\{X_1 > 1, 1 < X_2 \leq 4, Z \leq 2\}$
Λ_2	ℓ_2 : id = 2	$\{0 < X_1 \leq 1, X_2 \leq 3, Z \leq 5\}$	Λ_9	ℓ_9 : id = 9	$\{X_1 > 1, 1 < X_2 \leq 4, 2 < Z \leq 3\}$
Λ_3	ℓ_3 : id = 3	$\{X_1 \leq 1, X_2 > 3, Z \leq 5\}$	Λ_{10}	ℓ_{10} : id = 10	$\{X_1 > 1, X_2 \leq 4, Z > 5\}$
Λ_4	ℓ_4 : id = 4	$\{X_1 \leq 1, X_2 \leq 2, 5 < Z \leq 10\}$	Λ_{11}	ℓ_{11} : id = 11	$\{1 < X_1 \leq 5, X_2 > 4, Z \leq 7\}$
Λ_5	ℓ_5 : id = 5	$\{X_1 \leq 1, X_2 > 2, 5 < Z \leq 10\}$	Λ_{12}	ℓ_{12} : id = 12	$\{1 < X_1 \leq 5, X_2 > 4, Z > 7\}$
Λ_6	ℓ_6 : id = 6	$\{X_1 \leq 1, X_2 \in \mathbb{R}, Z > 10\}$	Λ_{13}	ℓ_{13} : id = 13	$\{X_1 > 5, X_2 > 4, Z \in \mathbb{R}\}$
Λ_7	ℓ_7 : id = 7	$\{X_1 > 1, X_2 \leq 1, Z \leq 3\}$			

Then, we would like to remove all nodes involving Z from $\mathcal{C}_{\mathbf{x},Z}$. In what follows, we present a step-by-step guide to the implementation of Algorithm 1, which is introduced in Section 3.2 of the main context. Figure 6 presents illustrative plots on the structure of the tree after each step. The removal process proceeds in a post-order.

- Step 1: The target node is $a = \{Z > 5\}$, removing the split of $\{Z \leq 10\}$ and $\{Z > 10\}$.

From the graph, we can see that one of a 's children (denoted by $c_1 = \{Z \leq 10\}$) is a sub-tree while the other child (denoted by $c_2 = \{Z > 10\}$) is a leaf node. We connect the children of c_1 to node a by setting $left(a) = left(c_1)$ and $right(a) = right(c_1)$, and then merge the node c_2 into the leaf nodes of c_1 , i.e. ℓ_4 and ℓ_5 in Figure 5. In specific, we modify the ids of ℓ_4 and ℓ_5 into $id(\ell_4) = id(\ell_4) \cup id(c_2) = 4 \cup 6$ and $id(\ell_5) = id(\ell_5) \cup id(c_2) = 5 \cup 6$. Figure (6a) plots the current tree structure after the removal.

- Step 2: The target node is $a = \{X_1 \leq 1\}$, removing the split of $\{Z \leq 5\}$ and $\{Z > 5\}$.

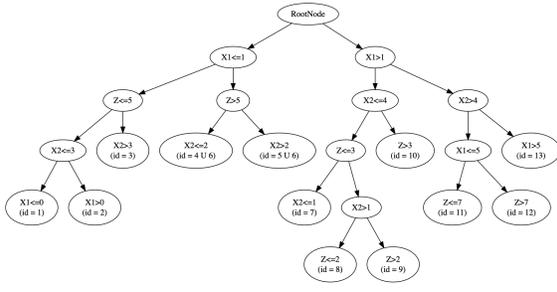
Both of a 's children are subtrees. From a high-level point of view, this process consists of two parts: (i) removing the splits associated with Z and merge the branches, followed by (ii) one breadth-first traversal to remove empty branches.

- Step 2(i): We connect the children of one subtree ($c_1 = \{Z \leq 5\}$) to a 's parent node by setting $left(a) = left(c_1)$ and $right(a) = right(c_1)$, and then connect the children of another subtree ($c_2 = \{Z > 5\}$) to every leaf node of c_1 , i.e., ℓ_1 , ℓ_2 and ℓ_3 in Figure (6a). More precisely, for $a' \in \{\ell_1, \ell_2, \ell_3\}$, we set $left(a') = left(c_2)$ and $right(a') = right(c_2)$, in the meantime, we update the ids of all the leaf nodes in c_2 (denoted by a'') by setting them to be $id(a'') = id(a') \cup id(a'')$. See Figure (6b1) for a graph illustration.
- Step 2(ii): We acknowledge that Step 2(i) may result in some infeasible regions because of contradictory conditions. We need one additional step to remove these

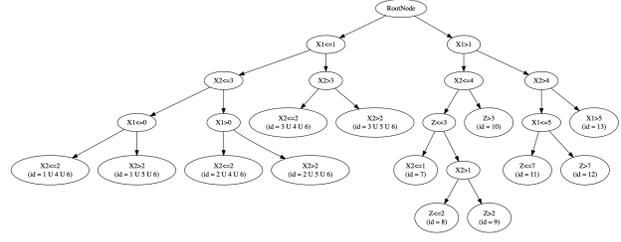
redundant nodes from the output. To this end, we conduct a level order traversal to remove nodes that correspond to empty regions. For example, in Figure (6b1), we observe that the node with $id = 3 \cup 4 \cup 6$ (denoted by c) corresponds to the path “ $X_1 \leq 1, X_2 > 3, X_2 \leq 2$ ”, which gives us an empty region as a result of contradictory conditions. Let a_2 denote its parent node, i.e. $a_2 = \{X_2 > 3\}$, and let c' denote the sibling of c , i.e., c' is the node with $id = 3 \cup 5 \cup 6$. Since c is empty, we remove the split of c and c' from the tree by setting $left(a_2) = left(c')$ and $right(a_2) = right(c')$. The fact that c' is a leaf node further makes a_2 becomes a new leaf node, whose id is set to as $id(a_2) = id(c') = 3 \cup 5 \cup 6$. The traversal will help delete infeasible branches so that the output produces non-empty and non-overlapping splits of \mathbf{X} . Figure (6b2) presents the final output.

- Step 3: The target node is $a = \{X_2 > 1\}$, removing the split of $\{Z \leq 2\}$ and $\{Z > 2\}$. Since both of a 's children are leaf nodes, we directly remove the two splits from the tree and convert a into a leaf node, whose id is set to be the union of the previous two leaf nodes, i.e., $id(a) = id(\{Z \leq 2\}) \cup id(\{Z > 2\}) = 8 \cup 9$. See Figure (6c).
- Step 4: The target node is $a = \{X_2 \leq 4\}$, removing the split of $\{Z \leq 3\}$ and $\{Z > 3\}$. This is similar to Step 1. One of a 's children is a subtree while the other is a leaf node, see Figure (6d) for details.
- Step 5: The target node is $a = \{X_1 \leq 5\}$, removing the split of $\{Z \leq 7\}$ and $\{Z > 7\}$. This is similar to Step 3 as both of a 's children are leaf nodes, see Figure (6e) for the structure after removing this split.

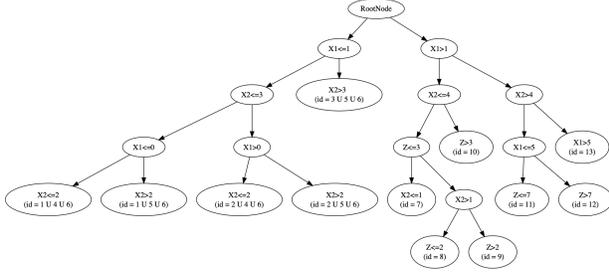
At this point, we have finished removing all the nodes associated with Z from the tree. The above output (denoted by $\mathcal{C}_{\mathbf{X}}$) produce a minimal partition with respect to \mathbf{X} . The resultant leaf cohorts, Π_1, \dots, Π_K with $K = 9$, are summarized in Table 4.



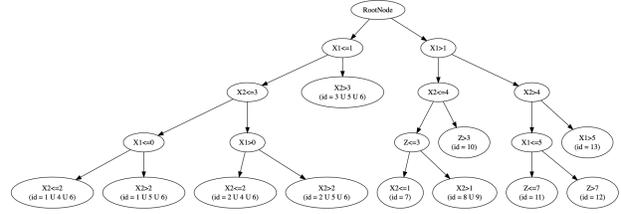
(a) Tree structure after removing the split of $\{Z \leq 10\}$ and $\{Z > 10\}$.



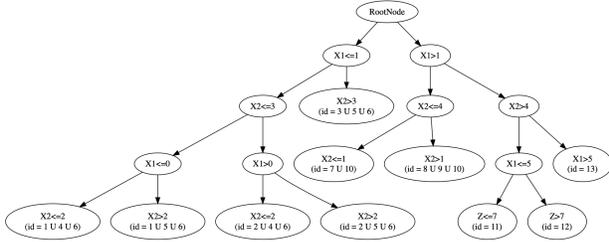
(b) Tree structure after removing the split of $\{Z \leq 5\}$ and $\{Z > 5\}$.



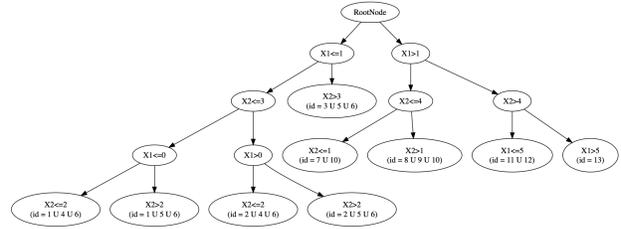
(b2) Tree structure after deleting infeasible branches



(c) Tree structure after removing the split of $\{Z \leq 2\}$ and $\{Z > 2\}$.



(d) Tree structure after removing the split of $\{Z \leq 3\}$ and $\{Z > 3\}$.



(e) Tree structure after removing the split of $\{Z \leq 7\}$ and $\{Z > 7\}$.

Figure 6: Plots of steps to remove Z from $\mathcal{C}_{\mathbf{X},Z}$ using Algorithm 1.

Note: The notation “id = 1 \cup 2” in the plots means that the samples in this merged node come from leaf nodes ℓ_1 and ℓ_2 in reference to the original tree $\mathcal{C}_{\mathbf{X},Z}$. Please be advised that this does not mean all the samples from the leaf nodes ℓ_1 and ℓ_2 are currently in this node with id = 1 \cup 2.

Similarly, by removing all nodes involving \mathbf{X} from $\mathcal{C}_{\mathbf{X},Z}$, we obtain a tree consisting of nodes related to Z , denoted by \mathcal{C}_Z . A minimal partition of $range(Z)$ can be extracted from the leaf-cohorts $\Gamma_1, \dots, \Gamma_L$. The results are summarized in Table 5.

Table 4: Leaf-cohorts Π_1, \dots, Π_K from the causal tree $\mathcal{C}_{\mathbf{X}}$

Leaf-cohorts	Leaf nodes	Cohort regions
Π_1	$a_{\mathbf{X},1}$: id = $1 \cup 4 \cup 6$	$\{X_1 \leq 0, X_2 \leq 2\}$
Π_2	$a_{\mathbf{X},2}$: id = $1 \cup 5 \cup 6$	$\{X_1 \leq 0, 2 < X_2 \leq 3\}$
Π_3	$a_{\mathbf{X},3}$: id = $2 \cup 4 \cup 6$	$\{0 < X_1 \leq 1, X_2 \leq 2\}$
Π_4	$a_{\mathbf{X},4}$: id = $2 \cup 5 \cup 6$	$\{0 < X_1 \leq 1, 2 < X_2 \leq 3\}$
Π_5	$a_{\mathbf{X},5}$: id = $3 \cup 5 \cup 6$	$\{X_1 \leq 1, X_2 > 3\}$
Π_6	$a_{\mathbf{X},6}$: id = $7 \cup 10$	$\{X_1 > 1, X_2 \leq 1\}$
Π_7	$a_{\mathbf{X},7}$: id = $8 \cup 9 \cup 10$	$\{X_1 > 1, 1 < X_2 \leq 4\}$
Π_8	$a_{\mathbf{X},8}$: id = $11 \cup 12$	$\{1 < X_1 \leq 5, X_2 > 4\}$
Π_9	$a_{\mathbf{X},9}$: id = 13	$\{X_1 > 5, X_2 > 4\}$

Table 5: Leaf-cohorts $\Gamma_1, \dots, \Gamma_L$ from the causal tree \mathcal{C}_Z

Leaf-cohorts	Cohort regions	Leaf-cohorts	Cohort regions
Γ_1	$\{Z \leq 2\}$	Γ_4	$\{5 < Z \leq 7\}$
Γ_2	$\{2 < Z \leq 3\}$	Γ_5	$\{7 < Z \leq 10\}$
Γ_3	$\{3 < Z \leq 5\}$	Γ_6	$\{Z > 10\}$

The next step is that, for each heterogeneous cohort Π_i , $i = 1 \dots K$, we would like to find the estimated causal effects for given treatment values, and report the optimal treatment associated with this heterogeneous cohort. For any $\mathbf{x} \in \text{range}(\mathbf{X})$ and $t \in \text{range}(T) \setminus \{0\}$, we can identify the unique cohort pair (Π_i, Γ_j) such that $\mathbf{x} \in \Pi_i$ and $t \in \Gamma_j$. We estimate the treatment effect as $\hat{\delta}(\mathbf{x}, t) = \hat{\delta}(\Pi_i, \Gamma_j)$ for any $\mathbf{x} \in \Pi_i$ and $t \in \Gamma_j$. What's more, for each pair of (Π_i, Γ_j) , we can find a unique leaf node ℓ_k in $\mathcal{C}_{\mathbf{X},Z}$, corresponding to a leaf-cohort Λ_k such that $\Pi_i \subseteq \Lambda_k(\mathbf{X})$ and $\Gamma_j \subseteq \Lambda_k(Z)$. The treatment effect can then be estimated as $\hat{\delta}(\Pi_i, \Gamma_j) = \hat{\delta}(\Lambda_k)$. Considering the leaf-cohorts in Tables 3 - 5, Table 6 provides a summary of relationships between (Π_i, Γ_j) and Λ_k .

One attractive feature of our proposed approach is that the task of estimating causal effects can be achieved simultaneously with the tree construction procedures from $\mathcal{C}_{\mathbf{X},Z}$ to $\mathcal{C}_{\mathbf{X}}$. The information of estimated causal effects has been encoded in the node name of leaf nodes in $\mathcal{C}_{\mathbf{X}}$, in which we keep track of the source leaf nodes in $\mathcal{C}_{\mathbf{X},Z}$ for the merged nodes in $\mathcal{C}_{\mathbf{X}}$. For example, for the leaf-cohort Π_1 , it is associated with the leaf node $a_{\mathbf{X},1}$ with $id = 1 \cup 4 \cup 6$. This means in order to estimate the causal effects in Π_1 with respect to

Table 6: Estimated causal effects: $\hat{\delta}(\Pi_i, \Gamma_j) = \hat{\delta}(\Lambda_k)$

Leaf-cohorts of $\mathcal{C}_{\mathbf{X}}$	Leaf-cohorts of \mathcal{C}_Z					
	Γ_1	Γ_2	Γ_3	Γ_4	Γ_5	Γ_6
Π_1	Λ_1	Λ_1	Λ_1	Λ_4	Λ_4	Λ_6
Π_2	Λ_1	Λ_1	Λ_1	Λ_5	Λ_5	Λ_6
Π_3	Λ_2	Λ_2	Λ_2	Λ_4	Λ_4	Λ_6
Π_4	Λ_2	Λ_2	Λ_2	Λ_5	Λ_5	Λ_6
Π_5	Λ_3	Λ_3	Λ_3	Λ_5	Λ_5	Λ_6
Π_6	Λ_7	Λ_7	Λ_{10}	Λ_{10}	Λ_{10}	Λ_{10}
Π_7	Λ_8	Λ_9	Λ_{10}	Λ_{10}	Λ_{10}	Λ_{10}
Π_8	Λ_{11}	Λ_{11}	Λ_{11}	Λ_{11}	Λ_{12}	Λ_{12}
Π_9	Λ_{13}	Λ_{13}	Λ_{13}	Λ_{13}	Λ_{13}	Λ_{13}

various treatment values, we only need to look into its three source codes, which corresponds to Λ_1 (Leaf ℓ_1), Λ_4 (Leaf ℓ_4) and Λ_6 (Leaf ℓ_6). The optimal treatment for cohort Π_1 is determined by

$$t^* = \begin{cases} 0 & \text{if } \max(\hat{\delta}(\Lambda_1), \hat{\delta}(\Lambda_4), \hat{\delta}(\Lambda_6)) \leq 0, \\ \arg \max_{t \in \text{range}(Z)} \hat{\delta}(\Pi_1, t) = \{t \in \Lambda_k \text{ for } k \text{ satisfying } \hat{\delta}(\Lambda_k) = \max(\hat{\delta}(\Lambda_1), \hat{\delta}(\Lambda_4), \hat{\delta}(\Lambda_6))\} & \text{otherwise.} \end{cases}$$

In summary, this observation enables us to estimate the treatment effects for one cohort in $\mathcal{C}_{\mathbf{X}}$ by only looking at its source nodes in $\mathcal{C}_{\mathbf{X},Z}$. This feature greatly facilitates the task of identifying the optimal treatment.

B Proof of Theorem 1

Under the assumptions given in Section 1, we have $\tau(z, \mathbf{x}) = E[Y | \mathbf{X} = \mathbf{x}, T = z] - E[Y | \mathbf{X} = \mathbf{x}, T = 0]$. It follows from the definition of W , Z_1 and Z that $E[Y | \mathbf{X} = \mathbf{x}, T = z] = E[Y | \mathbf{X} = \mathbf{x}, W = 1, Z_1 = z] = E[Y | \mathbf{X} = \mathbf{x}, W = 1, Z = z]$. Next, it follows from the independence of (Z, \mathbf{X}) and W that $E[Y | \mathbf{X} = \mathbf{x}, W = 1, Z = z] = E[Y(W = 1) | \mathbf{X} = \mathbf{x}, Z = z]$. Hence, $E[Y | \mathbf{X} = \mathbf{x}, T = z] = E[Y(W = 1) | \mathbf{X} = \mathbf{x}, Z = z]$.

Since $T = 0$ is equivalent to $W = 0$ (by definition), we have $E[Y | \mathbf{X} = \mathbf{x}, T = 0] = E[Y | \mathbf{X} = \mathbf{x}, W = 0] = E[Y(W = 0) | \mathbf{X} = \mathbf{x}]$. Finally, it follows from the independence of Z and

$(Y(W = 0), \mathbf{X})$ that $E[Y(W = 0) | \mathbf{X} = \mathbf{x}] = E[Y(W = 0) | \mathbf{X} = \mathbf{x}, Z = z]$. With the two, we have $E[Y | \mathbf{X} = \mathbf{x}, T = 0] = E[Y(W = 0) | \mathbf{X} = \mathbf{x}, Z = z]$. This completes the proof.

C Discussions on Simulation Designs

In this section, we introduce the intuition behind our simulation designs. We let

$$\eta(x_1, x_2) = -2 + \frac{2}{1 + \exp(-12(x_1 - 0.2))} \times \frac{2}{1 + \exp(-12(x_2 - 0.2))}$$

denote a function of heterogeneous features $\mathbf{X}_i = (X_{1i}, X_{2i})$. Figure (3) provides a graphical illustration of the values of $\eta(x_1, x_2)$ for (x_1, x_2) within a unit square. Blue and red depict the regions in which $\eta(x_1, x_2)$ takes positive and negative values, respectively. The heterogeneous nature of how the function $\eta(x_1, x_2)$ responds to features makes it the building block for designing our experiments. Besides $\eta(x_1, x_2)$, we include $\eta(x_1, 1 - x_2)$ to infuse more heterogeneity to the data. By defining the function of heterogeneous treatment effects, denoted by $f(x_1, x_2, t)$, as a function consisting of $\eta(x_1, x_2)$ and $\eta(x_1, 1 - x_2)$, different pairs of (x_1, x_2) will have various optimal treatments. We choose $f(x_1, x_2, t)$ to be the following functions under the different settings:

(i) Continuous:

$$\begin{aligned} f(x_1, x_2, t) = & 5\eta(x_1, x_2)\mathbb{1}_{\{t \in (0, 0.3]\}} - 5\eta(x_1, x_2)\mathbb{1}_{\{t \in (0.3, 0.5]\}} + \eta(x_1, 1 - x_2)\mathbb{1}_{\{t \in (0.5, 0.7]\}} \\ & - \eta(x_1, 1 - x_2)\mathbb{1}_{\{t \in (0.7, 1]\}}. \end{aligned}$$

(ii) Ordinal:

$$f(x_1, x_2, t) = 5\eta(x_1, x_2)\mathbb{1}_{\{t=1,2\}} - 5\eta(x_1, x_2)\mathbb{1}_{\{t=5\}} + \eta(x_1, 1 - x_2)\mathbb{1}_{\{t=3,4\}} - \eta(x_1, 1 - x_2)\mathbb{1}_{\{t=6\}}.$$

(iii) Categorical:

$$f(x_1, x_2, t) = 5\eta(x_1, x_2)\mathbb{1}_{\{t=a\}} - 5\eta(x_1, x_2)\mathbb{1}_{\{t=b\}} + \eta(x_1, 1 - x_2)\mathbb{1}_{\{t=c\}} - \eta(x_1, 1 - x_2)\mathbb{1}_{\{t=d\}}.$$

Using the categorical case as an example, Table 7 summarizes the values of heterogeneous treatment effects $f(x_1, x_2, t)$ with four different pairs of heterogeneous features (x_1, x_2) . From the table we can see that for (x_1, x_2) at different locations in the unit square, it responds distinctively with respect to different treatments. As a result, the optimal treatment varies across various pairs of heterogeneous features.

$$f(x_1, x_2, t) = \begin{cases} 5\eta(x_1, x_2), & t = a. \\ -5\eta(x_1, x_2), & t = b. \\ \eta(x_1, 1 - x_2), & t = c. \\ -\eta(x_1, 1 - x_2), & t = d. \\ 0, & t = 0. \end{cases}$$

Table 7: Values of $f(x_1, x_2, t)$ with different pairs of (x_1, x_2)

(x_1, x_2)	Treatment t				optimal treatment for (x_1, x_2)
	a	b	c	d	
(0.8, 0.8)	9.970	-9.970	-0.001	0.001	a
(0.2, 0.8)	-0.007	0.007	-1.000	1.000	d
(0.8, 0.2)	-0.007	0.007	1.994	-1.994	c
(0.2, 0.2)	-5.000	5.000	-0.001	0.001	b