

Causal Modeling with Neural Networks and Individualized Treatment Effect Estimation

Master Thesis in Biostatistics (STA495)

by

Mike Krähenbühl

Matriculation number: 18-652-149

supervised by

Prof. Dr. Beate Sick, University of Zurich & ZHAW

Prof. Dr. Oliver Dürr, HTWG Konstanz

Zurich, July 2025

**Causal Modeling with Neural
Networks
and
Individualized Treatment Effect
Estimation**

Mike Krähenbühl

Version July 22, 2025

Contents

Preface	iii
Abstract	v
1 Introduction	1
1.1 Motivation	1
1.2 Key concepts of causal inference	2
1.3 Goals and contributions	3
2 Methods	5
2.1 TRAM-DAGs	5
2.2 Individualized Treatment Effects (ITEs)	12
2.3 Software	14
3 Experiments	15
3.1 Experiment 1: TRAM-DAG (simulation)	15
3.2 Experiment 2: ITE on International Stroke Trial (IST)	20
3.3 Experiment 3: ITE model robustness in RCTs (simulation)	24
3.4 Experiment 4: ITE estimation with TRAM-DAGs (simulation)	32
4 Discussion	55
4.1 Findings	55
4.2 Limitations	56
4.3 Conclusion	56
Bibliography	57
5 Appendix	61
5.1 Interpretation of linear coefficients	61
5.2 Bernstein polynomial for continuous outcomes	61

5.3	Negative log-likelihood	63
5.4	Encoding of discrete variables	64
5.5	Scaling of continuous variables	65
5.6	Modeling interactions with complex shift (CS)	66
5.7	Experiment 2: Calibration plots	68
5.8	Experiment 3: Standard Random Forest for ITE Estimation	71
5.9	Experiment 3: Calibration plots	72
5.10	Declaration of tools and services used	73

Preface

This thesis marks the final part of my Master of Science in Biostatistics at the University of Zurich. I wanted to work on a topic where I could apply my interest and deepen my knowledge in machine learning, especially in relation to causal questions.

The TRAM-DAG framework ([Sick and Dürr, 2025](#)), developed by my supervisors Prof. Dr. Beate Sick and Prof. Dr. Oliver Dürr, provided a perfect opportunity to do so. Our initial aim was to apply it to real-world data and potentially include semi-structured data. However, due to some surprising findings by [Chen *et al.* \(2025\)](#), our focus shifted towards the increasingly important topic of individualized treatment effect (ITE) estimation. Towards the end, we then bridged ITE estimation with the TRAM-DAG framework.

I want to thank my supervisors and all the people I had the chance to work and study with, as well as everyone who supported me on this journey.

Mike Krähenbühl
July 2025

Abstract

This thesis explores the use of TRAM-DAGs, a flexible neural network-based framework for estimating complex causal relationships in known directed acyclic graphs (DAGs). TRAM-DAGs allow sampling from observational, interventional, and counterfactual distributions when the full DAG is known. We applied TRAM-DAGs to both simulated data and a real-world randomized controlled trial, with a focus on individualized treatment effect (ITE) estimation.

We show how TRAM-DAGs can be used with continuous and ordinal or categorical predictor variables, investigate how variable scaling affects interpretability, and demonstrate how interactions between variables can be modeled. A key part of this work involved applying different causal machine learning models – including TRAM-DAGs, logistic regression, and random forests – to estimate ITEs on the International Stroke Trial (IST) dataset. In line with findings by [Chen *et al.* \(2025\)](#), none of the models produced ITE estimates that generalized to the test data.

To explore possible reasons for this poor performance, we conducted simulation experiments under varying conditions. These revealed that, while models must be well calibrated and overfitting should be avoided, this alone may not suffice to guarantee valid ITE estimates. Weak treatment-covariate interaction effects and especially unmeasured effect modifiers were found to be critical challenges. When such variables are unobserved, the ignorability assumption alone may not ensure unbiased estimation – an issue also highlighted by [Vegetabile \(2021\)](#). These factors may help explain the limited model performance observed in the IST dataset.

We also applied TRAM-DAGs in randomized and confounded simulation settings with relatively complex DAGs and found that, when the full DAG was observed and interaction effects were present, TRAM-DAGs accurately recovered causal relationships and provided unbiased ITE estimates.

While promising, our work has limitations. The simulation scenarios may not fully capture real-world complexity, and evaluating ITE estimates on real data remains challenging since ground truth is unknown. The neural network-based TRAM-DAGs require training time and rely on modeling assumptions – e.g., regarding the scale of conditional effects – when parameter interpretability is desired.

TRAM-DAGs offer a customizable modeling framework that enables the specification of both flexibility and interpretability, making them suitable for real-world causal inference tasks. Future work could apply TRAM-DAGs to more diverse datasets, including semi-structured data, and further investigate ITE estimation in the presence of unmeasured effect modifiers.

Keywords: TRAM-DAGs, neural causal model, individualized treatment effect, structural causal model, counterfactuals, transformation model, observational data, heterogeneous treatment effect, conditional average treatment effect

Chapter 1

Introduction

1.1 Motivation

The most important questions in research are mostly not associational, but causal ([Pearl, 2009a](#)). They concern the effects of interventions – such as the impact of a treatment – or seek explanations for observed outcomes, such as identifying which disease caused certain symptoms. They also include hypothetical scenarios; for example: what would the gross domestic product have been if interest rates had increased by 25 instead of 75 basis points? Answering such questions requires causal reasoning and demands an understanding of the underlying data-generating process. Purely associational approaches are typically not sufficient to draw valid causal conclusions.

The gold standard for estimating the causal effect of an intervention on an outcome is the randomized controlled trial (RCT) ([Hariton and Locascio, 2018](#)). In this prospective study design, participants are randomly assigned to either the treatment or control group. Randomization aims to eliminate the influence of potential confounding variables, ensuring that treatment groups are balanced with respect to baseline characteristics. This allows for an unbiased estimation of the causal effect. Despite their strengths, RCTs have several limitations. They are often expensive and time-consuming to plan and execute. Moreover, the results may not generalize well to the population of interest, as individuals who volunteer or are accepted for trials are not always representative of the target group. Another central limitation of RCTs is that in many scenarios they simply cannot be conducted due to ethical or practical reasons. For example, an RCT is only ethical in the case of clinical equipoise, which means that there is uncertainty about the superiority of one of the two treatment arms ([Freedman, 1987](#)). It is not acceptable to treat one group with the assumed inferior treatment. The same is true for obviously harmful interventions, such as smoking or drinking alcohol. In these cases, it is not possible to conduct an RCT to estimate the causal effect of smoking on lung cancer.

For these reasons, much of research aims to make causal inference from observational data, using non-experimental or quasi-experimental designs. Unlike RCTs, these settings do not involve randomization to treatment, which introduces challenges due to confounding. Methods for causal inference from observational data aim to correctly control for such confounders to enable valid causal conclusions.

An application where causal inference is of particular importance is the estimation of personalized treatment effects. RCTs typically estimate an average treatment effect (ATE) on a sample, which is the difference in mean outcomes between the treatment and the control groups ([Nichols, 2007](#)). However, individual patients may respond differently to the treatment, depending on their unique characteristics. In personalized medicine, the estimate of treatment effects at the individual level is referred to as the individualized treatment effect (ITE) or conditional average treatment effect (CATE), while in business and marketing contexts, the term uplift modeling is often used ([Gutierrez and Gérardy, 2017](#); [Zhao and Harinen, 2020](#)). Such estimates are critical in settings where treatment responses vary significantly between individuals. For clinical

decision-making, tailoring therapies to individual characteristics can lead to more effective and efficient care. The importance of estimating individual-level effects is not limited to medicine. It is also of high interest in marketing, where campaigns can be precisely targeted to maximize impact and minimize adverse responses. Consider, for instance, the decision of whether to send a push notification (treatment) to a customer. Some customers might be persuadables, who will respond positively only if treated. Others, in contrast, might have responded positively without the intervention but are negatively affected by receiving it – for example, a customer who is reminded of a forgotten subscription and, as a result, decides to cancel it. In this context, identifying persuadables is valuable, while treating the latter may be counterproductive. This illustrates the need to understand treatment effects at a granular level to guide individualized decisions. Various methods have been proposed to estimate individualized treatment effects, yet this task remains challenging. The fundamental problem is that only one of the two potential outcomes can ever be observed for any given individual (Holland, 1986), making the estimation of treatment effects inherently more difficult than standard predictive modeling.

These two motivations – causal inference from observational data and the estimation of individualized treatment effects — form the focus of this thesis. We build on existing methods to address both challenges. The next sections provide an overview of key concepts and methods, and lay out the specific research questions we aim to answer.

1.2 Key concepts of causal inference

Questions in causal inference are typically classified into one of the three levels of Pearl’s hierarchy of causation (Pearl, 2009b). Level 1 corresponds to observational queries, expressed as conditional probabilities $P(Y | X)$, which can be answered directly from the joint distribution $P(Y \cap X)$. Level 2 involves interventional queries, such as $P(Y | \text{do}(X = \alpha))$, which describe the probability when intervening on a variable X and setting it to a particular value. Level 3 addresses counterfactual reasoning. These are hypothetical what-if questions that require reasoning about outcomes under alternative realities. For example, if a patient received a treatment and died, the factual outcome is death under the received treatment. The counterfactual would be the outcome that would have occurred had the patient received a different treatment.

Causal relationships can be represented by a directed acyclic graph (DAG), where the variables, or nodes, are connected by directed edges, which represent causal dependencies.

While DAGs capture the structure of these dependencies, structural causal models (SCMs) extend this representation by explicitly modeling the functional relationships between variables. A set of structural equations of the form $X_i = f_i(\text{pa}(X_i), Z_i)$, $i = 1, \dots, n$ defines an SCM (Pearl, 2009b). Here, $\text{pa}(X_i)$ denotes the direct causal parents of X_i , and Z_i is an exogenous noise variable. These exogenous variables capture latent factors that influence X_i but are not explicitly modeled. By convention, the Z_i are assumed to be mutually independent. This assumption implies that the DAG includes all involved variables and causal connectivities. Any dependence between variables must be due to a direct or indirect path in the graph; otherwise, dependencies would show up in the noise terms, violating their independence (Pearl, 2009b).

Each function f_i – which may be nonlinear – defines how the value of X_i is generated from its parents and the corresponding noise term. A source node X_j without any parents is modeled as $X_j = f_j(Z_j)$. Once all structural equations and noise variables are specified, the model is fully deterministic in the sense that each variable is a fixed function of its parents and its own exogenous noise. The randomness in the system arises entirely from these independent noise terms. This functional representation makes it possible to compute interventional distributions and evaluate counterfactual outcomes. These aspects are discussed in detail in Section 2.1.4.

In this thesis, we do not focus on discovering the underlying causal graph. Such a structure may be obtained through structure learning algorithms (see e.g., Zheng *et al.*, 2018) or determined from expert knowledge. Instead, we assume the graph is known and fully observed, and

concentrate on estimating the functional form of the relationships between variables – that is, the structural equations that define the SCM.

Various approaches exist for estimating the functions f_i that constitute an SCM, depending on the assumptions made about the data and the model class.

A simple approach to modeling the structural equations is linear regression, which assumes Gaussian error terms Z_i and linear functional forms f_i . Classical statistical methods of this kind are typically well-defined, computationally efficient, and offer interpretable parameters. However, they rely on strong assumptions about the underlying data-generating mechanism – such as linearity and homoscedasticity – which often do not hold in practice. Violations of these assumptions can lead to biased or misleading results.

Alternatively, more flexible approaches based on tree-based ensemble models or neural networks have gained popularity for estimating structural equations. These models are capable of approximating complex, nonlinear relationships and capturing complicated interactions between variables with minimal bias. Their flexibility, however, often comes at the cost of reduced interpretability and, in some cases, limited applicability to non-continuous or mixed data types. [Poinsot et al. \(2024\)](#) provided an overview of deep structural causal models and their use in counterfactual inference.

The TRAM-DAG framework proposed by [Sick and Dürr \(2025\)](#) builds a bridge between these classical and neural-network-based modeling approaches, by combining interpretable transformation models with the flexibility of neural networks. At its core, the structural equations are modeled using transformation models ([Hothorn et al., 2014](#)), a flexible class of distributional regression methods. These models were subsequently extended to deep transformation models (Deep TRAMs) by [Sick et al. \(2021\)](#), enabling the use of neural networks to estimate the parameters of the transformation function in a flexible and customizable way. In the TRAM-DAG framework, these deep TRAMs are applied according to a known causal graph, allowing the model to be fitted to observational data and used to answer causal queries across all three levels of Pearl’s hierarchy. The framework is introduced in more detail in Section 2.1.3.

1.3 Goals and contributions

The first part of this thesis focuses on the TRAM-DAG framework. In particular, we apply the model to DAGs of varying complexity, extend it from continuous to ordinal and categorical predictors (see Section 5.4), examine how variable scaling affects the interpretation of coefficients (see Section 5.5), and make use of different neural network configurations (e.g., activation functions, batch normalization, dropout). Most analyses are conducted on simulated data, where the underlying data-generating process is known (see Sections 3.1 and 3.4), but the model is also applied to real-world data to demonstrate its practical utility (see Section 3.2).

The second focus of the thesis is the estimation of individualized treatment effects (ITEs). Recent work by [Chen et al. \(2025\)](#) showed that most causal machine learning models trained on RCT data failed to generalize when evaluated out of sample. We replicate part of their study by applying various models they used, as well as TRAM-DAGs, to the same data and analyze whether we reach similar conclusions (see Section 3.2). We further investigate why ITE estimation may fail in such settings and under which conditions reliable estimates can be obtained (see Section 3.3). In addition, we demonstrate that TRAM-DAGs can be effectively used to estimate ITEs also in complex, non-randomized observational settings, provided the causal graph is known and fully observed (see Section 3.4). In doing so, we explore the potential of TRAM-DAGs as a framework for answering causal questions across different levels of Pearl’s hierarchy.

Formally, we aim to answer the following research questions in this thesis:

- How can TRAM-DAGs be applied to estimate causal relationships across different sce-

narios, and subsequently be used to sample from observational, interventional, and counterfactual distributions? Using synthetic data, we demonstrate how the model can be applied (Section 3.1) and explain how to handle ordinal and categorical predictors (Section 5.4), how variable scaling affects the interpretation of coefficients (Section 5.5), and how interactions between variables can be modeled (Section 5.6).

- Do we reach similar conclusions as [Chen *et al.* \(2025\)](#) when applying their models, as well as TRAM-DAGs, to the International Stroke Trial (IST) dataset for ITE estimation? Specifically, we examine whether ITEs estimated on the training data fail to generalize to the test data (Section 3.2).
- What factors contribute to the failure of ITE estimation in causal machine learning models, and under which conditions can reliable estimates be obtained? We investigate this by applying various models across simulated scenarios that differ in observed variables and treatment effect size (Section 3.3).
- Can TRAM-DAGs provide unbiased estimates of ITEs when the causal graph is fully observed? We evaluate this in a complex simulation setting with a continuous outcome, comparing both randomized and confounded designs (Section 3.4).

With this work, we aim to contribute to the important and evolving field of causal inference in observational settings and to the challenging task of estimating individualized treatment effects.

Chapter 2

Methods

This chapter introduces the methodological foundations for the experiments used in this thesis. Section 2.1 presents the concept and functionality of TRAM-DAGs, along with the necessary theoretical background. Section 2.2 provides the framework for estimating individualized treatment effects. These methods form the basis for the experiments and analyses that follow in Chapter 3.

2.1 TRAM-DAGs

The goal of TRAM-DAGs is to estimate the structural equations of a given DAG in a flexible and, if desired, interpretable way. This enables the sampling of observational and interventional distributions, as well as make counterfactual statements. The approach requires data and a known causal graph describing the underlying structure. It is assumed that there are no hidden confounders. TRAM-DAGs estimate, for each variable X_i , a transformation function $Z_i = h_i(X_i \mid \text{pa}(X_i))$, where Z_i denotes the noise variable and $\text{pa}(X_i)$ are the causal parents of X_i . Crucially, for continuous variables, this relationship can be inverted to yield the structural equation $X_i = h_i^{-1}(Z_i \mid \text{pa}(X_i))$. The monotonically increasing transformation functions h_i represent the conditional distribution $X_i \mid \text{pa}(X_i)$ on a latent scale Z_i . They are based on the framework of transformation models as introduced by [Hothorn et al. \(2014\)](#) and were extended to deep TRAMs by [Sick et al. \(2021\)](#). The following summarizes the key ideas of these models, which form the building blocks of TRAM-DAGs.

2.1.1 Transformation Models

Transformation models are a flexible class of distributional regression models that can be applied to various data types. A linear transformation model is given by

$$F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z(h(y \mid \mathbf{x})) = F_Z(h_I(y) - \mathbf{x}^\top \boldsymbol{\beta}) \quad (2.1)$$

where $F_{Y|\mathbf{X}=\mathbf{x}}(y)$ is the conditional cumulative distribution function (CDF) of the outcome variable Y given the predictors $\mathbf{X} = \mathbf{x}$. The transformation function $h(y \mid \mathbf{x})$ maps the outcome variable y onto the latent scale of variable Z , and F_Z is the CDF of Z , the so-called inverse-link function that maps $h(y \mid \mathbf{x})$ to probabilities. In its simplest form, as shown in Equation 2.1, the transformation function can be split into an intercept part $h_I(y)$ and a linear shift component $\mathbf{x}^\top \boldsymbol{\beta}$, where \mathbf{x} are the predictors and $\boldsymbol{\beta}$ the corresponding coefficients.

If the latent distribution Z is chosen to be the standard logistic, each coefficient β_i can be interpreted as log-odds ratio: an increase of one unit in the predictor x_i , while holding other predictors unchanged, increases the log-odds (latent scale) of the outcome Y by β_i . After applying the inverse-link function F_Z , this can induce a non-linear change to the conditional

distribution of Y on the original scale. Further discussion on the choice of latent distribution and the interpretation of coefficients is provided in Appendix 5.1.

For a continuous outcome Y , the intercept $h_I(y)$ is represented by a Bernstein polynomial, which is a flexible and monotonically increasing function:

$$h_I(y) = \frac{1}{M+1} \sum_{k=0}^M \vartheta_k B_{k,M}(y), \quad (2.2)$$

where $\vartheta_k, k = 0, \dots, M$ are the coefficients and $B_{k,M}(y)$ are the corresponding Bernstein basis polynomials. The coefficients ϑ_k are constrained to be monotonically increasing to ensure that the transformation function $h_I(y)$ is strictly increasing. More details on the technical implementation of the Bernstein polynomial (Equation 2.2) in the context of deep-TRAMs are provided in the Appendix 5.2

For a discrete outcome Y , the intercept h_I is represented by a set of cut-points ϑ_k , which define the thresholds separating the different levels of the outcome. For example, a binary outcome requires one cut-point, while an ordinal outcome with K levels requires $K - 1$ cut-points. The transformation model is given by:

$$P(Y \leq y_k \mid \mathbf{X} = \mathbf{x}) = F_Z(\vartheta_k + \mathbf{x}^\top \boldsymbol{\beta}), \quad k = 1, 2, \dots, K - 1. \quad (2.3)$$

A visual representation for a continuous and discrete (ordinal) outcome is shown in Figure 2.1.

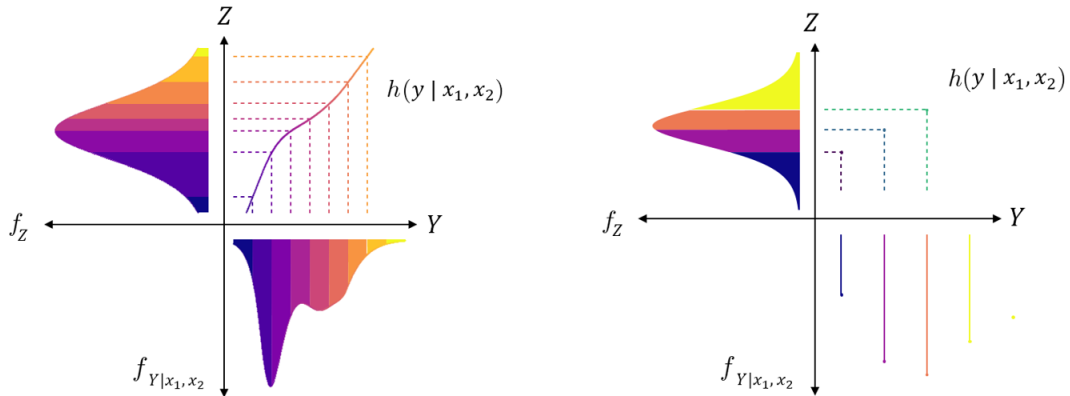


Figure 2.1: Left: Example of a transformation model for a continuous outcome Y (Equation 2.1) with a smooth transformation function. Right: Example of a transformation model for an ordinal outcome Y (Equation 2.3) with 4 levels. The transformation function consists of cut-points separating the levels of the outcome. In both cases the latent distribution Z is the standard logistic and the predictors \mathbf{x} induce a linear (vertical) shift of the transformation function. Plots adapted from visualizations similar to those used in Sick and Dürr (2025).

The parameters $\boldsymbol{\beta}$ and $\boldsymbol{\vartheta}$ are estimated by minimizing the negative log-likelihood (NLL), defined as:

$$\text{NLL} = -\frac{1}{n} \sum_{i=1}^n \log \left(f_{Y|\mathbf{X}=\mathbf{x}_i}(y_i) \right), \quad (2.4)$$

where $f_{Y|\mathbf{X}=\mathbf{x}_i}(y_i)$ is the conditional density function of the outcome Y given the predictors \mathbf{x} of the i -th observation under the current parameterization. A full derivation is provided in Appendix 5.3.

Throughout this thesis, these transformation models form the foundation for estimating conditional distributions of variables within the TRAM-DAG framework. Unless stated otherwise, the latent distribution F_Z is chosen to be the standard logistic, resulting in a logistic transformation model.

2.1.2 Deep TRAMs

The transformation models introduced in Section 2.1.1 were extended into deep TRAMs using a modular neural network architecture (Sick *et al.*, 2021). The goal is to obtain a parameterized transformation function of the form

$$h(y | \mathbf{x}_L, \mathbf{x}_C) = h_I(y) + \mathbf{x}_L^\top \boldsymbol{\beta} + f(\mathbf{x}_C), \quad (2.5)$$

where $h_I(y)$ is the intercept (simple or complex), \mathbf{x}_L are predictors with a linear effect on the transformation function, and \mathbf{x}_C are those with a potentially complex, non-linear influence. The parameters for each component of the model – intercept, linear shift, and complex shift – are obtained by a neural network module. The user can assign predictors to these components depending on the assumed structure of the data and the desired flexibility. Figure 2.2 illustrates an example of a transformation function (Equation 2.5) with these three components (SI-LS-CS). For simplicity, it shows the simple intercept (SI) case where the intercept does not depend on predictors.

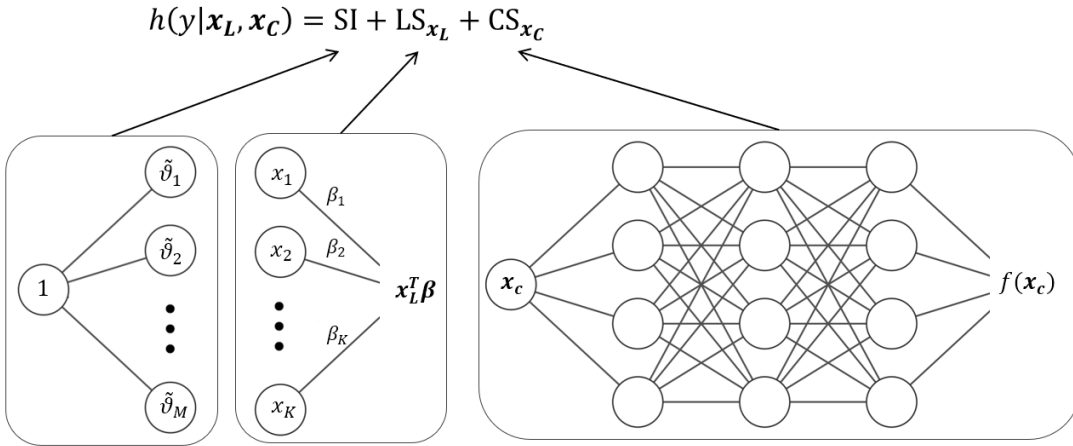


Figure 2.2: Modular deep transformation model (deep TRAM). The transformation function $h(y | \mathbf{x})$ (Equation 2.5) is constructed from the outputs of three neural networks: intercept (SI/CI), linear shift (LS), and complex shift (CS). The intercept module (left) shows the simple intercept (SI) case and could be extended to a complex intercept (CI) by feeding predictors into the module and adding hidden layers with nonlinear activation functions.

Intercept: The intercept defines the baseline shape of the transformation function when $\mathbf{x}_L^\top \boldsymbol{\beta} = 0$ and $f(\mathbf{x}_C) = 0$. For continuous outcomes, it is modeled using a smooth Bernstein polynomial (Equation 2.2), and for discrete outcomes via cut-points. The preliminary parameters $\hat{\vartheta}_k$ are the output nodes of a neural network, and subsequently transformed into monotonically increasing parameters ϑ_k . A simple intercept (SI) assumes no dependency on predictors and uses only a constant input. A complex intercept (CI), by contrast, allows the parameters to vary as a function of selected predictors by feeding \mathbf{x} into a neural network with hidden layers, producing predictor-dependent parameters $\vartheta_k(\mathbf{x})$. This allows the baseline distribution to adapt flexibly to different covariate settings. Details on the computation of the Bernstein polynomial (Equation 2.2) are provided in Appendix 5.2.

Linear shift: Predictors with an assumed linear influence on the transformation function are modeled via the linear shift (LS) component. Here, a neural network without hidden layers and without bias nodes receives \mathbf{x}_L as input and returns a linear combination $\mathbf{x}_L^\top \boldsymbol{\beta}$ as output. This results in a vertical, linear shift of the transformation function. The parameters $\boldsymbol{\beta}$ are interpretable coefficients, and in the case of a logistic transformation model, represent log-odds ratios. See Appendix 5.1 for more details.

Complex shift: Non-linear dependencies between predictors and the transformation function can be modeled by the complex shift (CS) component. The corresponding predictors \mathbf{x}_C are input into a deep neural network (with one or more hidden layers and non-linear activation functions), yielding a scalar output $f(\mathbf{x}_C)$. This allows modeling of non-linear predictor effects, including interactions between predictors if multiple predictors are input into the neural network module. See Appendix 5.6 for an example of modeling an interaction with a complex shift for the purpose of estimating individualized treatment effects (ITEs).

Level of complexity: A key advantage of the deep TRAM architecture is that users can specify the role of each predictor: linear (LS), complex (CS), or influencing the shape of the transformation function (CI). For example, Kook *et al.* (2022) demonstrated the use of deep transformation models with flexible components for ordinal outcomes and multi-modal data. Herzog *et al.* (2023) later applied this framework in the medical domain, predicting ordinal stroke outcomes by combining structured tabular data with image features (via a CNN). This design enables the integration of different data modalities into a single model.

The resulting distribution function of deep TRAMs is invariant to the choice of the inverse-link function F_Z (latent distribution) in an unconditional (Hothorn *et al.*, 2018) or fully flexible (CI) setting. However, once restrictions are introduced on the influence of the predictors – such as LS or CS – the model assumes a fixed scale of dependence. The choice of latent distribution may depend on (i) assumptions about the data-generating process, (ii) the conventional and widely used interpretation scale for parameters (e.g., log-odds ratios or log-hazard ratios), and (iii) a data-driven criterion such as selecting the distribution that minimizes the negative log-likelihood (NLL) on the observed data.

Parameter estimation: The weights of the neural networks are learned by minimizing the NLL (Equation 2.4) of the conditional model. Starting from random weight initialization, the networks are trained using the Adam optimizer (Kingma and Ba, 2015). Outputs of the modules are assembled to compute the NLL, and backpropagation (Rumelhart *et al.*, 1986) is used to adjust weights iteratively. Deep learning techniques such as dropout (Srivastava *et al.*, 2014), early stopping (Prechelt, 2012), and batch normalization (Ioffe and Szegedy, 2015) can be applied to prevent overfitting and improve generalization. Nonlinear activation functions (e.g., ReLU Glorot *et al.*, 2011 or sigmoid Rumelhart *et al.*, 1986) are used in hidden layers to capture complex relationships. Schmidhuber (2015) provides an extensive historical overview of deep learning methods and their development.

2.1.3 TRAM-DAGs: Deep TRAMS applied in a causal setting

In TRAM-DAGs, deep transformation models (see Section 2.1.2) are applied within a causal setting. We assume a pre-specified DAG that represents the causal dependencies among variables. Figure 2.3 illustrates the basic idea using a simple DAG with three variables and no hidden confounders. Each node’s conditional distribution is modeled using a deep TRAM, given its parent variables in the DAG. The assumed influence of the parent variables must be specified as CI, LS or CS. In this example, X_1 is a continuous source node (i.e., without parents) and its transformation function consists only of a simple intercept (SI). X_2 is also continuous and depends on X_1 by a linear shift (LS). X_3 is an ordinal variable with four levels; its transformation function is influenced linearly by X_1 (LS) and non-linearly by X_2 (CS). The cut-points

$h(x_{3,k} | x_1, x_2)$ represent the cumulative log-odds of levels $k = 1, 2, 3$ of X_3 . The probability of the highest category ($k = 4$) is the complement of the cumulative probability of the first three.

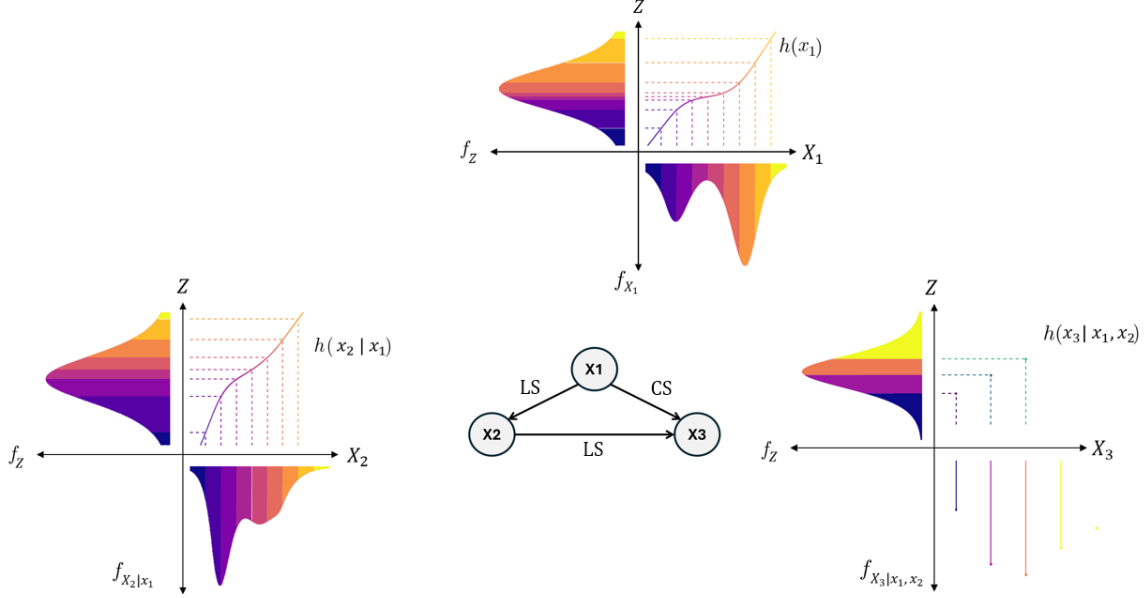


Figure 2.3: Example of a TRAM-DAG with three variables: X_1 , X_2 , and X_3 . Each variable’s distribution is modeled conditionally on its parents using deep TRAMs. Arrows indicate the causal dependencies. Plots adapted from visualizations similar to those used in Sick and Dürr (2025).

The DAG in Figure 2.3, along with its assumed dependencies, can be represented by a meta-adjacency matrix (Equation 2.6), where rows indicate sources and columns indicate targets of causal influence:

$$\mathbf{MA} = \begin{bmatrix} 0 & \text{LS} & \text{LS} \\ 0 & 0 & \text{CS} \\ 0 & 0 & 0 \end{bmatrix} \quad (2.6)$$

Based on this DAG, each variable’s conditional distribution is modeled using a deep TRAM, enabling generative sampling and causal inference.

Constructing the modular neural network: As described in Section 2.1.2, the transformation functions are implemented using a modular neural network. The network takes as input the observed variables \mathbf{x} and the meta-adjacency matrix (Equation 2.6), which encodes the assumed causal structure. This matrix acts as a binary mask that restricts information flow to follow the directed edges of the graph, ensuring that only causal parent-to-child relationships are modeled. As a result, spurious dependencies such as $X_2 \rightarrow X_1$ are prevented. This masking approach is inspired by the Masked Autoregressive Flow (MAF) framework (Papamakarios *et al.*, 2017), where connectivity within the network is masked to control which inputs each output can depend on.

Before training the model, the input data must be appropriately pre-processed. Discrete variables with few categories are dummy encoded, and continuous variables should be scaled prior to input. Further discussion of encoding and scaling is provided in Appendix 5.4 and 5.5.

Once the data are preprocessed and the structure is specified, the architecture of the neural networks for the complex shift or complex intercept must be defined – including choices such as depth, width, activation functions, and whether to use dropout or batch normalization. These

decisions depend on the assumed complexity of the effects and the need to regularize against overfitting.

Each node's transformation function (Equation 2.5) is assembled from the outputs of its three modular components: the intercept (SI or CI), linear shift (LS), and complex shift (CS). Model training is performed by minimizing the NLL (Equation 2.4), optimizing all parameters at the same time. The estimated parameters β in the linear shift are interpretable as log-odds ratios corresponding to a one-unit increase in the respective parent variable, holding all other parent variables unchanged.

2.1.4 Sampling from TRAM-DAGs

Once a TRAM-DAG is fitted on data, it can be used to sample from the observational or interventional distribution, or to perform counterfactual queries. For continuous outcomes, the structural equations $X_i = f(Z_i, \text{pa}(X_i))$ are represented by the inverse of the conditional transformation functions, i.e., $X_i = h^{-1}(Z_i | \text{pa}(X_i))$, since the transformation function maps from observed values to the latent scale: $Z_i = h(X_i | \text{pa}(X_i))$.

Observational sampling: The sampling process from the observational distribution is described in Algorithm 1 and illustrated in Figure 2.4. In each iteration, one complete sample of all variables in the DAG is generated – i.e., a sample from the joint observational distribution.

Algorithm 1 Generate a complete sample from the observational distribution

Given: A fitted TRAM-DAG with structural equations $X_i = h^{-1}(Z_i | \text{pa}(X_i))$
for each node X_i in topological order **do**
 Sample latent value $z_i \sim F_{Z_i}$ ▷ e.g., `rlogis()` in R
 if X_i is continuous **then**
 Solve $h(x_i | \text{pa}(x_i)) - z_i = 0$ for x_i ▷ numerical root-finding
 else if X_i is discrete **then**
 Find category $x_i = \max(\{0\} \cup \{x : z_i > h(x | \text{pa}(x_i))\}) + 1$
 end if
end for

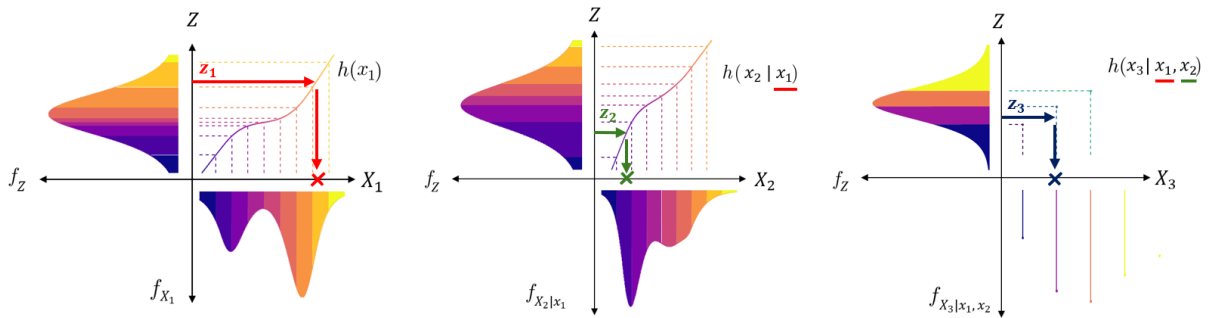


Figure 2.4: One sampling iteration for the three variables from the estimated transformation functions $h(x_i | \text{pa}(x_i))$. The latent values z_i are sampled from the standard logistic distribution. The values x_i are determined by applying the inverse of the transformation function for continuous variables or by finding the corresponding category for the ordinal variable. See Algorithm 1 for details. Plots adapted from visualizations similar to those used in Sick and Dürr (2025).

Interventional sampling: To sample from an interventional distribution, we apply the *do*-operator as described by Pearl (1995). The *do*-operator fixes a variable at a specific value, thereby

removing causal dependencies on its (former) parents. This allows us to simulate interventions by deleting all edges pointing into the intervened variable, while sampling the remaining variables as described earlier in Algorithm 1.

Counterfactual queries: In a counterfactual query, we are interested in what the value of a variable X_i would have been, had another variable X_j taken a different value than actually observed. Pearl (2009b) describes a three-step procedure to answer such queries. In short, let \mathbf{x} denote the observed values (a sample) of all variables in a DAG and let \mathbf{z} denote the corresponding latent variables inferred via the transformation functions: $z_k = h_k(x_k \mid \text{pa}(x_k))$. Then, in a counterfactual query where X_j is set to a new value α (intervention), we use the same latent values \mathbf{z} , but apply the transformation functions in a post-interventional DAG with $X_j := \alpha$ to compute the counterfactual outcomes for the remaining variables. The counterfactual estimation procedure is outlined in Algorithm 2 and visualized in Figure 2.5. Note that this visualization is simplified to illustrate only one parent variable (X_j); in practice, the set of parents $\text{pa}(X_i)$ may include additional predictors.

Algorithm 2 Answer a single counterfactual query

Given: A TRAM-DAG model with estimated structural equations $X_k = h_k^{-1}(Z_k \mid \text{pa}(X_k))$

Input: Observed sample \mathbf{x} , intervention $X_j := \alpha$

Step 1 (Abduction): For each observed variable x_k , compute the corresponding latent value $z_k = h_k(x_k \mid \text{pa}(x_k))$

Step 2 (Action): Modify the DAG by setting $X_j := \alpha$ (intervention)

Step 3 (Prediction): Using the counterfactual DAG and the intervened value $X_j = \alpha$, go along the causal order and determine the counterfactual values for all descendants X_k of X_j . This is done by evaluating

$$x_k = h_k^{-1}(z_k \mid \text{pa}^*(x_k)),$$

where z_k is the latent value obtained in Step 1, and $\text{pa}^*(x_k)$ are the (possibly updated) parent values after the intervention.

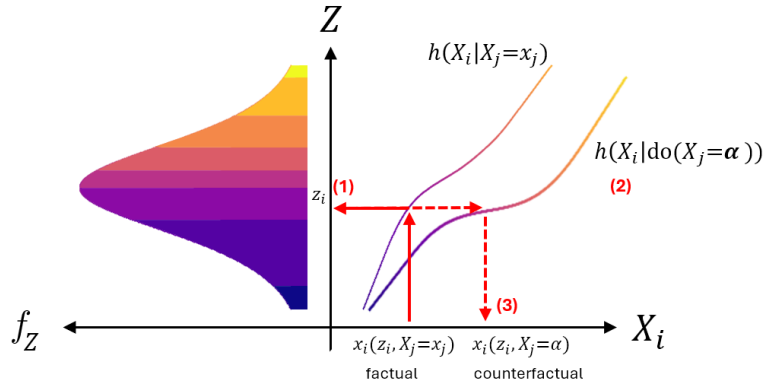


Figure 2.5: Illustration of a counterfactual query for variable X_i , had its parent X_j taken a different value, following the three-step procedure: (1) abduction of the latent value z_i from the observed outcome, (2) intervention by setting $X_j = \alpha$, (3) prediction of the counterfactual outcome using the same z_i and the modified transformation function. This simplified example assumes a single parent variable; in practice, $\text{pa}(X_i)$ may include multiple predictors. Plot adapted from visualizations similar to those used in Sick and Dürr (2025).

While interventional distributions (e.g., $P(X_i = x_i \mid \text{do}(X_j = \alpha), z_i)$) are well-defined in

both continuous and discrete settings, determining the counterfactual realizations is generally only possible for continuous variables. In the discrete case, no unique latent value z_k can be recovered from an observed category, as multiple z_k values may map to the same outcome.

2.2 Individualized Treatment Effects (ITEs)

While the homogeneous treatment effect refers to the part of the effect that is equal for all patients, the heterogeneous treatment effect (HTE) describes the variation in treatment effects across individuals or subgroups. One reason for treatment effect heterogeneity is the presence of interaction variables (Hoogland *et al.*, 2021), also referred to as effect modifiers (Christensen *et al.*, 2021), which influence how the treatment effect varies depending on individual characteristics. However, even in the absence of explicit interaction terms, treatment effects can still exhibit heterogeneity in certain model classes, such as logistic regression (Hoogland *et al.*, 2021), or in other subclasses of transformation models. This phenomenon is further explored in Experiment 4 and will be discussed in detail in Section 3.4.3.2.

In personalized medicine and informed decision-making, individual-level treatment effects are also of interest besides population averages. The individual treatment effect compares the potential outcomes for an individual i under two different treatment options: the potential outcome $Y_i(1)$ under treatment, and $Y_i(0)$ under control, as defined by Rubin’s potential outcomes framework (Rubin, 2005). The individual treatment effect is typically defined as the difference between the two potential outcomes:

$$Y_i(1) - Y_i(0). \quad (2.7)$$

However, for any individual, only one of these two potential outcomes can be observed – the other one remains counterfactual. This is known as the fundamental problem of causal inference (Holland, 1986), which implies that the individual treatment effect (Equation 2.7) is inherently unobservable. Estimating causal effects from observed data requires untestable assumptions to ensure identifiability, in contrast to standard predictive modeling.

Assumptions for identifiability: A central assumption is consistency, which requires that the observed outcome equals the potential outcome under the treatment actually received: $Y = Y(1)$ if $T = 1$, and $Y = Y(0)$ if $T = 0$. The Stable Unit Treatment Value Assumption (SUTVA) assumes no interference between units and that treatments are well-defined (Rubin, 1980). The most critical assumption is ignorability (or unconfoundedness), which states that, conditional on observed covariates \mathbf{X} , the treatment assignment is independent of the potential outcomes (Rosenbaum and Rubin, 1983):

$$(Y(1), Y(0)) \perp T \mid \mathbf{X}. \quad (2.8)$$

Equation 2.8 implies that there are no unmeasured confounders affecting both treatment assignment and outcomes. In an RCT, this condition is typically satisfied by design. In contrast, for observational data, the covariate set \mathbf{X} must contain a valid adjustment set to appropriately control for confounding. Additionally, care must be taken to avoid adjusting for collider variables, as this can introduce selection bias (see Elwert and Winship (2014) and Hernan and Robins (2025)). Finally, the positivity assumption requires that every individual has a non-zero probability of receiving each treatment level:

$$0 < P(T = 1 \mid \mathbf{X} = \mathbf{x}) < 1 \quad \text{for all } \mathbf{x}. \quad (2.9)$$

While the individual treatment effect is inherently unobservable, its expectation conditional on covariates can be identified under the assumptions above. This quantity is referred to as the conditional average treatment effect (CATE), or individualized treatment effect (ITE). In this

thesis, we use the term ITE to denote this estimand. For an individual with covariates $\mathbf{X} = \mathbf{x}_i$, the ITE can be identified as:

$$\begin{aligned}
 \text{ITE}(\mathbf{x}_i) &= \mathbb{E}[Y_i(1) - Y_i(0) \mid \mathbf{X} = \mathbf{x}_i] \\
 &= \mathbb{E}[Y_i(1) \mid \mathbf{X} = \mathbf{x}_i] - \mathbb{E}[Y_i(0) \mid \mathbf{X} = \mathbf{x}_i] \\
 &= \mathbb{E}[Y_i(1) \mid T = 1, \mathbf{X} = \mathbf{x}_i] - \mathbb{E}[Y_i(0) \mid T = 0, \mathbf{X} = \mathbf{x}_i] \quad (\text{by ignorability}) \\
 &= \mathbb{E}[Y_i \mid T = 1, \mathbf{X} = \mathbf{x}_i] - \mathbb{E}[Y_i \mid T = 0, \mathbf{X} = \mathbf{x}_i] \quad (\text{by consistency})
 \end{aligned} \tag{2.10}$$

In the case of a binary outcome Y , where the expected value corresponds to the probability that $Y = 1$, the ITE for an individual with covariates $\mathbf{X} = \mathbf{x}_i$, under the identifiability assumptions, is given by:

$$\text{ITE}(\mathbf{x}_i) = P(Y_i = 1 \mid T = 1, \mathbf{X} = \mathbf{x}_i) - P(Y_i = 1 \mid T = 0, \mathbf{X} = \mathbf{x}_i). \tag{2.11}$$

If $Y = 1$ indicates recovery, a positive ITE means that the individual is expected to benefit from the treatment compared to control. A negative ITE suggests a lower probability of recovery under treatment, while an ITE of zero indicates no expected difference between treatment and control for that individual.

ITE estimation: In practice, the ITE (Equation 2.10) is estimated by fitting a model to the observed data and then predicting the outcome under both treatment conditions. For binary outcomes, this requires a model that predicts the probability of $Y = 1$ under treatment and control (i.e., the expected value of a binary outcome; see Equation 2.11). For continuous outcomes, a model is required that directly predicts the expected value of each potential outcome. Alternatively, one may compute the difference in the predicted medians of the potential outcomes, corresponding to the quantile treatment effect (QTE) at the median (Equation 3.2). This is illustrated in the simulation example in Experiment 4 (see Section 3.4). An overview of the models used for ITE estimation is provided below.

Models for ITE estimation: While many other approaches for ITE estimation exist, this thesis focuses on two widely used types of metalearners: the T-learner and the S-learner (Künzel *et al.*, 2019). These methods are model-agnostic, meaning that any predictive algorithm can be used as a base model. The T-learner fits two separate models: one on the treated group and one on the control group. If the treatment was assigned randomly, confounding should not be a problem. The ITE for an individual is then estimated as the difference between the predicted outcomes from these two models (as in Equation 2.10). In contrast, the S-learner fits a single model on all individuals, incorporating the treatment indicator as an additional covariate to distinguish between treatment conditions. Therefore, the S-learner requires a model that is complex enough to capture interaction effects between treatment and covariates.

The experiments in this thesis were conducted using both S- and T-learners. For binary outcomes, we used models based on logistic regression, logistic lasso regression, random forests, and TRAM-DAGs to predict the probabilities of $Y = 1$ for the potential outcomes (see Experiments 2 and 3; Sections 3.2 and 3.3). For continuous outcomes, we used TRAM-DAGs, which—when applied as S-learners—can flexibly model treatment effect heterogeneity when specified with complex shift or complex intercept terms. For simplicity, the ITEs for continuous outcomes were computed as the difference in the predicted medians of the potential outcomes (Equation 3.2), although estimation based on expected values would also have been possible, as discussed in the ITE estimation procedure in Experiment 4 (see Section 3.4.2).

Accurate ITE estimation requires models that generalize well, which in turn means avoiding overfitting. Since ITEs are computed as differences between two outcome predictions, even small

errors can accumulate. Hoogland *et al.* (2021) emphasize the need to balance model complexity with data availability to avoid overfitting. They also highlight that estimating risk differences is inherently more challenging than predicting single outcomes. Furthermore, calibration is crucial when ITEs are derived from predicted potential outcomes (Hoogland *et al.*, 2024). While neural networks may be accurate, they often produce poorly calibrated probabilities (Guo *et al.*, 2017), which can mislead treatment effect estimation. Similarly, conventional models like logistic regression can overfit in small samples or high-dimensional settings, leading to extreme predictions. Penalization methods such as lasso shrink coefficients to improve generalization (Riley *et al.*, 2021), though studies have shown that they can perform inconsistently in low-sample scenarios (Calster *et al.*, 2020). These factors are important to consider when estimating ITEs.

When using random forests as S-learners, care must be taken to ensure that the treatment variable is included in the splits. Proper tuning of hyperparameters like tree depth and variable selection (e.g., `mtry`) may be necessary to avoid overfitting. In summary, model choice for ITE estimation should consider the complexity of treatment-outcome relationships, sample size, and calibration needs. While complex models like TRAM-DAGs and random forests can capture complicated heterogeneity, simpler models may offer more robust performance in smaller datasets or where treatment-covariate interactions are straightforward.

Validation of ITE estimation: Hoogland *et al.* (2024) examined a variety of evaluation methods for assessing discrimination and calibration in ITE estimation. They emphasize that accurate prediction of potential outcomes should be considered a prerequisite for reliable ITE modeling. Furthermore, as in standard prediction modeling, proper validation should ideally be performed on independent (out-of-sample) data. In this thesis, although various quantitative evaluation metrics exist, we primarily relied on visual validation tools. One such approach is the ITE-ATE plot (see, e.g., Figure 3.9), where individuals are grouped based on their predicted ITEs. Within each ITE-subgroup, the empirical ATE is calculated using observed outcomes. Depending on the context, the empirical ATE is calculated on the same scale as the predicted ITEs – for example, as a risk difference for binary outcomes or an absolute difference for continuous outcomes. If the ITE estimates are well-calibrated, the observed ATEs should align with the predicted values, resulting in a calibration curve that lies along the identity line. Additionally, in simulation experiments where the true ITEs are known, estimation accuracy was assessed by directly comparing predicted and true effects. Furthermore, model calibration in terms of predictive performance $P(Y \mid \mathbf{X} = \mathbf{x})$ was evaluated using calibration plots. As an additional check, the average of the estimated ITEs across the population should equal the estimated ATE, i.e., $\mathbb{E}[\text{ITE}(X)] = \text{ATE}$. Large deviations from this condition may indicate systematic bias. However, satisfying this equality does not guarantee that the estimated ITEs are accurate at the individual level. All evaluations were conducted both on the training data and on a separate hold-out test set.

2.3 Software

All code used in this thesis is available on GitHub: https://github.com/mikekr97/MA_Mike.

All analyses were conducted in R (R Core Team, 2024) (version 4.4.2) using RStudio. The packages `keras` (Chollet *et al.*, 2017) (version 2.15.0), `tensorflow` (Allaire and Tang, 2025) (version 2.16.0), and `reticulate` (Ushey *et al.*, 2025) (version 1.40.0) were used to build and train neural networks through Python’s TensorFlow backend. These tools allowed for the use of deep learning methods directly within the R environment.

Chapter 3

Experiments

In this chapter, we present four experiments, each designed to address one of the research questions outlined in Section 1.3. Specifically:

- **Experiment 1** serves as a proof of concept to demonstrate how TRAM-DAGs can recover causal relationships in a known DAG and subsequently answer causal queries from all three levels of Pearl’s causal hierarchy. This is illustrated using simulated data (see Section 3.1).
- **Experiment 2** recreates the ITE estimation study by [Chen *et al.* \(2025\)](#) on the RCT data from the International Stroke Trial (IST). We apply three causal ML models, including TRAM-DAGs, to assess whether we reach similar conclusions – namely, that the estimated ITEs do not generalize to independent test data (see Section 3.2).
- **Experiment 3** investigates, via simulation, which factors may cause ITE estimation to fail, as observed in the IST dataset (see Section 3.3). Specifically, we simulate an RCT with a binary outcome under three scenarios: (3.1) an ideal case with full observability and strong heterogeneity (Section 3.3.3.1), (3.2) a case with an unobserved effect modifier (Section 3.3.3.2), and (3.3) a case with weak heterogeneity (Section 3.3.3.3).
- **Experiment 4** explores whether TRAM-DAGs can provide unbiased ITE estimates in complex settings, assuming the full DAG is observed. We conduct a simulation study with a continuous outcome under both randomized and confounded conditions, using a complex DAG structure. The simulations are conducted under three scenarios: (4.1) an ideal case with direct treatment effect and interaction effects (Section 3.4.3.1), (4.2) a case including a direct effect but without interaction effects (Section 3.4.3.2), and (4.3) a case with interaction effects but no direct effect (Section 3.4.3.3). If TRAM-DAGs yield unbiased estimates, this supports their utility for ITE estimation in realistic scenarios (see Section 3.4).

Together, these experiments aim to provide insights into the application of TRAM-DAGs and related causal ML methods for individualized treatment effect estimation, as well as general limitations encountered in practice.

3.1 Experiment 1: TRAM-DAG (simulation)

3.1.1 Motivation

This experiment demonstrates the application of TRAM-DAGs on a synthetic dataset, using the illustrative DAG previously shown in Figure 2.3. The objective is to show how TRAM-DAGs can learn causal relationships from observational data, assuming a known DAG. After fitting

the model to the joint distribution generated by the underlying causal structure, it is used to sample from observational, interventional, and counterfactual distributions.

The controlled simulation setting allows us to interpret the learned model components in detail and evaluate the model’s ability to recover both linear and nonlinear causal relationships.

3.1.2 Setup

We visualized the model fitting in terms of the training loss and subsequently showed and interpreted the learned components of the transformation functions (Equation 2.5), such as intercepts, linear and complex shifts. Finally, we drew samples from the estimated distributions to obtain observational and interventional distributions. We also conducted counterfactual queries on the learned model.

Data-generating process: We simulated a dataset with three variables, X_1 , X_2 , and X_3 , following the structure of the DAG and its associated meta-adjacency matrix shown in Figure 3.1. The matrix describes the functional dependencies between variables, where LS indicates a linear shift and CS a complex shift. Rows represent the source of the effect, and columns the target.

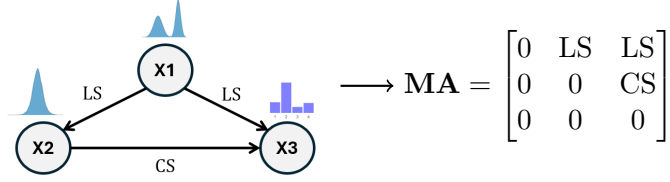


Figure 3.1: Causal graph (left) and meta-adjacency matrix (right) for Experiment 1 (Section 3.1). The transformation function of X_2 depends on X_1 via a linear shift (LS). The transformation function of X_3 depends on X_1 via a linear shift (LS) and on X_2 via a complex shift (CS).

The variable X_1 is continuous and bimodally distributed, and acts as a source node in the DAG, i.e., it is not influenced by any other variable:

$$X_1 = \begin{cases} \mathcal{N}(0.25, 0.1^2) & \text{with probability 0.5,} \\ \mathcal{N}(0.73, 0.05^2) & \text{with probability 0.5} \end{cases}$$

The second variable, X_2 , is continuous and linearly dependent on X_1 on the log-odds scale, with a true coefficient of $\beta_{12} = 2$. Its transformation function is

$$h(X_2 | X_1) = h_I(X_2) + \beta_{12}X_1,$$

where the baseline transformation (i.e., intercept) of X_2 is $h_I(X_2) = 5X_2$.

The third variable, X_3 , is ordinal and depends on both X_1 (LS) and X_2 (CS). We define the complex shift induced by X_2 as $f(X_2) = 0.5 \cdot \exp(X_2)$, and specify the linear shift parameter for X_1 as $\beta_{13} = 0.2$. The transformation function for category k of the ordinal variable X_3 with 4 levels (K) is thus defined by

$$h(X_{3,k} | X_1, X_2) = \vartheta_k + \beta_{13}X_1 + f(X_2),$$

with cut-points $\vartheta_k \in \{-2, 0.42, 1.02\}$ defining the thresholds of the ordinal variable. We generated samples for X_2 and X_3 as described in Section 2.1.4, by first sampling a latent value from the standard logistic distribution and then determining the corresponding observation using the transformation function.

This simulation allows us to assess whether the TRAM-DAG model can correctly recover the functional forms of the conditional dependencies and the associated parameters (linear and complex).

Model: Given the meta-adjacency matrix and the simulated observations, we construct a modular neural network based on the TRAM-DAG framework. The complex shift from X_2 to X_3 is modeled using a neural network with 4 hidden layers and 2 nodes per layer, as illustrated in Figure 3.2. A total of 20,000 samples are generated according to the defined DGP to fit the model. We train the model for 400 epochs using the Adam optimizer (Kingma and Ba, 2015) with a learning rate of 0.005.

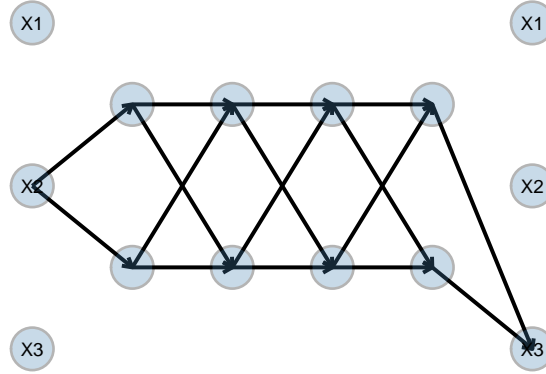


Figure 3.2: Neural network architecture for the complex shift on X_3 from X_2 in Experiment 1 (Section 3.1). The complex shift is modeled by a neural network with 4 hidden layers of shape (2, 2, 2, 2), using non-linear activation functions (sigmoid).

Model evaluation: We compare the estimated intercepts, learned coefficients, and the complex shift to the true values used in the DGP. We also compare the sampled observational and interventional distributions to the true distributions. For the counterfactual queries, we show the estimated counterfactual values for X_2 under an intervention on X_1 at a specific value, and compare these to the true counterfactual outcomes.

3.1.3 Results

We evaluate whether TRAM-DAGs can recover the structural equations and distributions used in the data-generating process described in Section 3.1.2. Below, we present the training process and inspect the estimated parameters, distributions, and counterfactual predictions.

Figure 3.3 shows the loss and the estimated parameters for the linear shifts over epochs during training. The loss was minimized during training and the estimated parameters β_{12} and β_{13} converged to the true values used in the DGP. The linear shift parameters are the interpretable part of the model (log-odds ratios). From the fitted model, we generated samples from the observational distribution, as shown in Figure 3.4. The TRAM-DAG can recover the observational distribution as its samples align with the data that was used to fit the model. Then we drew samples from the interventional distribution, where $X_2 = 1$ is fixed, as shown in Figure 3.5. Fixing X_2 leads to a distributional change in X_3 , which was also captured by the model. The TRAM-DAG learns the linear shifts (β_{12} , β_{13}) and the complex shift $f(X_2)$, which are shown in Figure 3.6. Figure 3.7 presents the intercepts learned for each of the nodes. For comparison, we added the estimated intercept functions from the Continuous Outcome Logistic

Regression (`Colr()`) function from the `tram` R-package (Hothorn *et al.*, 2018) for X_1 and X_2 , and the true values used in the DGP for the ordinal variable X_3 (three cut-points for the four levels). Since the transformation functions for X_1 and X_2 contain no complex terms, they match the default form used in `Colr()`. Finally, Figure 3.8 shows the counterfactuals for X_2 estimated by the TRAM-DAG for varying values of X_1 . The counterfactuals are the predicted values of X_2 had X_1 taken other values instead of the initially observed one.

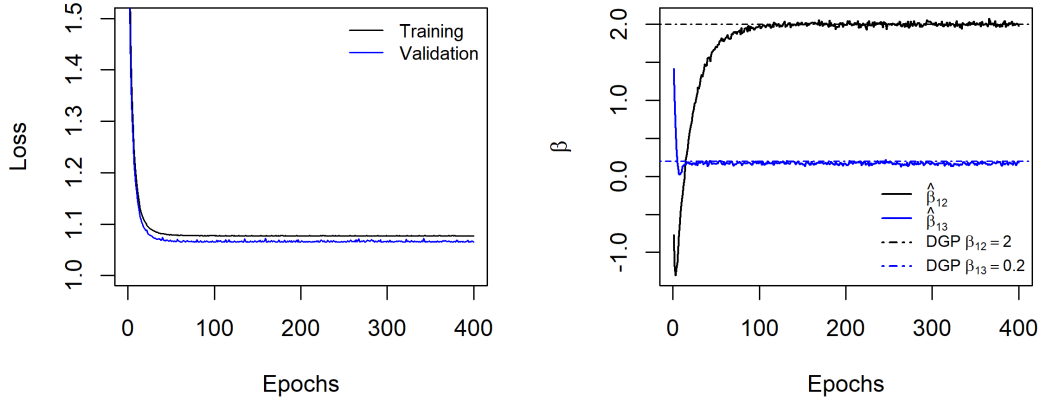


Figure 3.3: TRAM-DAG model fitting over 400 epochs for Experiment 1. Left: loss functions on the training and validation sets; Right: estimated parameters (betas) for the linear shift components over epochs. The estimates converge to the true values used in the DGP.

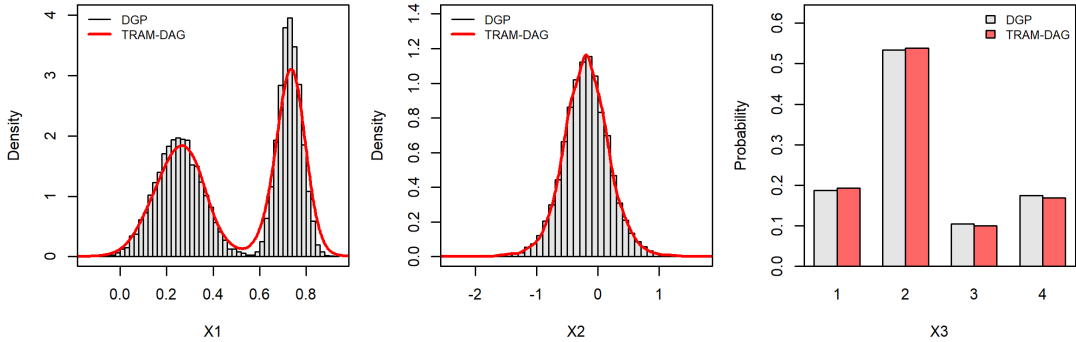


Figure 3.4: Samples generated by the TRAM-DAG from the learned observational distribution, compared to the true observations from the DGP.

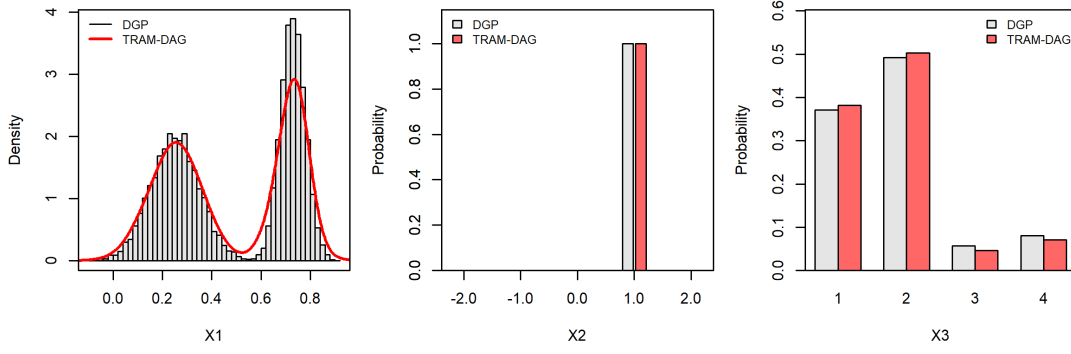


Figure 3.5: Samples generated by the TRAM-DAG compared to the true observations from the interventional distribution of the DGP, where $X_2 = 1$ is fixed. According to the DAG, this intervention induces a distributional change in X_3 .

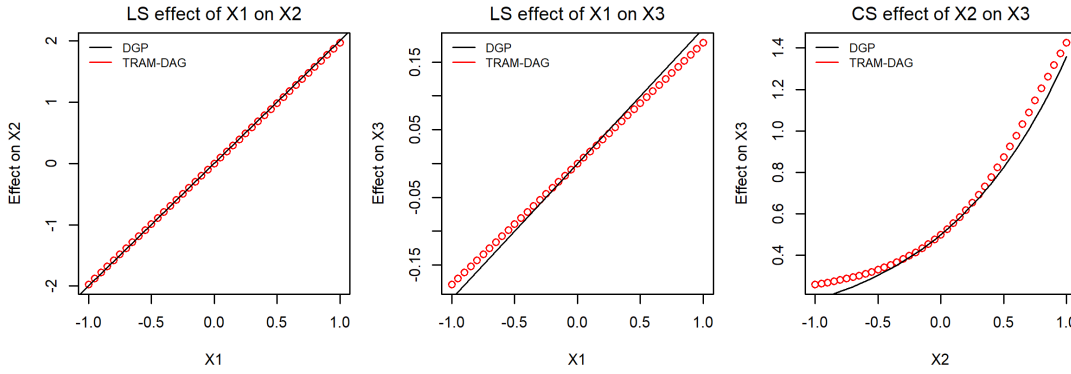


Figure 3.6: Linear and complex shifts learned by the TRAM-DAG. Left: $LS(X_1)$ on X_2 ; Middle: $LS(X_1)$ on X_3 ; Right: $CS(X_2)$ on X_3 . For visualization, we subtracted $\delta_0 = CS(0) - f(0)$ from the estimated complex shift $CS(X_2)$ to align it with the true shift function $f(X_2)$ from the DGP.

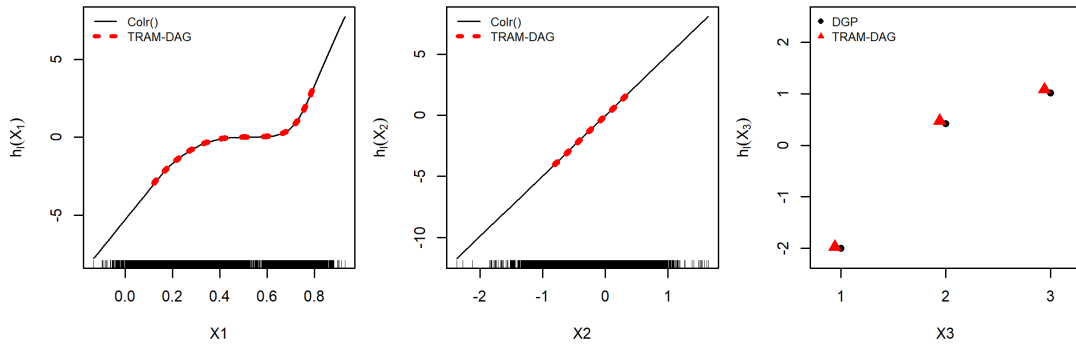


Figure 3.7: Intercepts learned for each of the nodes, along with the estimates from the `Colr()` function for the continuous variables and the true values from the DGP for the ordinal variable X_3 . Left: Smooth baseline transformation function for continuous X_1 ; Middle: Smooth baseline transformation function for continuous X_2 ; Right: Cut-points as the baseline transformation function for ordinal X_3 . For the last plot, we added $\delta_0 = CS(0) - f(0)$ to the estimated cut-offs to make them comparable to the true parameters from the DGP.

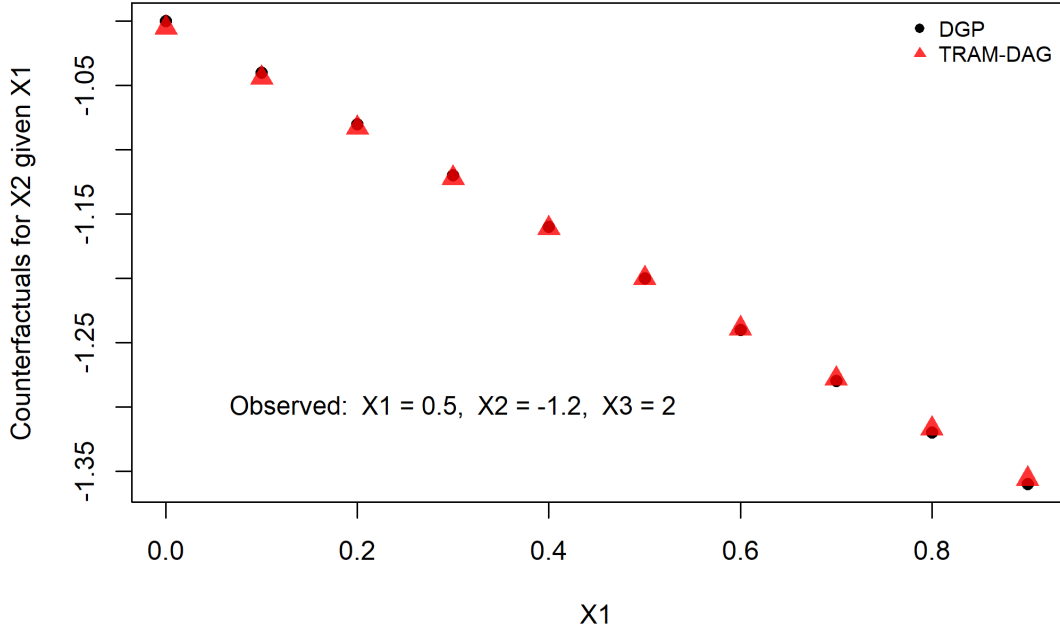


Figure 3.8: Counterfactuals for X_2 estimated with the TRAM-DAG for varying values of X_1 . We assumed the observed values $X_1 = 0.5$, $X_2 = -1.2$, and $X_3 = 2$, and determined the counterfactual values of X_2 had X_1 taken different values instead of the actually observed one. This illustrates how the model estimates alternative outcomes under hypothetical interventions on X_1 .

3.1.4 Discussion of Experiment 1

The results show that the TRAM-DAG framework can accurately recover the underlying causal relationships from data, including both linear and nonlinear (complex) relationships. This enabled the model to generate observational and interventional distributions (Figures 3.4 and 3.5) that closely matched the true ones, and produce valid counterfactual predictions (Figure 3.8).

This experiment provides a simple proof of concept for the flexibility and generative capability of TRAM-DAGs. By incorporating both interpretable (e.g., linear shifts) and flexible (e.g., complex shift) components, the model is able to capture a range of causal mechanisms – provided that the true DAG is known and data is generated accordingly.

3.2 Experiment 2: ITE on International Stroke Trial (IST)

3.2.1 Motivation

Chen *et al.* (2025) evaluated multiple causal ML methods on the International Stroke Trial (IST), to estimate the individualized treatment effects (ITEs). They demonstrated that none of the applied ML methods generalized well, as performance on the test data differed significantly from the training data on the chosen evaluation metrics. In this experiment, we replicate the analysis on the same data by applying three causal ML methods for ITE estimation, to investigate whether we obtain similar results as the authors.

3.2.2 Setup

Data: The International Stroke Trial was a large, randomized controlled trial conducted in the 1990s to assess the efficacy and safety of early antithrombotic treatment in patients with acute

ischemic stroke (International Stroke Trial Collaborative Group, 1997). Using a 2x2 factorial design, 19,435 patients across 36 countries were randomized within 48 hours of symptom onset to receive aspirin, subcutaneous heparin, both, or neither. Patients allocated to aspirin (300 mg daily for 14 days) had a 6-month death or dependency rate of 62.2%, compared to 63.5% in the control group not receiving aspirin, corresponding to a statistically significant absolute risk reduction after adjustment for baseline prognosis (1.4%, $p = 0.03$). The authors stated that there was no interaction between aspirin and heparin in the main outcomes. In this thesis, we focus exclusively on the aspirin vs. no aspirin comparison and the outcome of death or dependency at 6 months after stroke.

The dataset used in this experiment was made publicly available by Sandercock *et al.* (2011) and contains individual-level data, including baseline covariates assessed at randomization, treatment allocation, and 6-month outcomes, with a follow-up rate of 99%.

We used the same data pre-processing steps as Chen *et al.* (2025) to ensure comparability of results. 5.9% of individuals had incomplete data and were removed from the dataset. We used 2/3 of the data for fitting the models and 1/3 as a hold out test set. The final dataset included 21 baseline variables recorded at randomization: aspirin allocation (treatment), age, delay between stroke and randomization (in hours), systolic blood pressure, sex, CT performed before randomization, visible infarct on CT, atrial fibrillation, aspirin use within 3 days prior to randomization, and presence or absence of neurological deficits (including face, arm/hand, leg/foot deficits, dysphasia, hemianopia, visuospatial disorder, brainstem or cerebellar signs, and other neurological deficits), as well as consciousness level, stroke subtype, and geographical region. The outcome variable was death or dependence at 6 months.

Models for ITE estimation: The aim is to estimate the ITE (Equation 2.11) based on baseline characteristics. As a benchmark, we apply a T-learner logistic regression (following Chen *et al.* (2025), using the `stats` package). As a more complex model, we apply a T-learner tuned random forest (using the `comets` package (Kook, 2024)), which tunes the number of variables considered for splitting at each node (`mtry`) and the maximum tree depth (`max.depth`) using out-of-bag error, with 500 trees. Additionally, we apply an S-learner TRAM-DAG. For the random forest and TRAM-DAG based methods, we additionally scale numerical and dummy encode categorical covariates prior to model training. The transformation function of the outcome is modelled by a complex intercept $h(Y | T, \mathbf{X}) = CI(T, \mathbf{X})$, with 4 hidden layers of shape (20, 10, 10, 2). This architecture allows for interaction between the treatment and covariates. Furthermore, batch normalization, ReLU activation, and dropout (0.1) are applied to prevent overfitting and stabilize learning. A validation set comprising 20% of the training data is used to select the model with the lowest out-of-sample negative log-likelihood (Equation 2.4), while the test set remains untouched for final evaluation.

Since the IST is a randomized controlled trial, the full potential of TRAM-DAGs – designed primarily for use in observational settings – is not required here, as only the outcome needs to be modeled as a function of baseline patient characteristics. However, applying TRAM-DAGs in this context still allows us to assess its predictive performance and ability to flexibly model interactions between variables.

Model evaluation: For validation, since the ground truth is not known, we first rely on calibration plots to assess the general prediction power for the probabilities. Second, we predict the potential outcomes with the trained models to estimate the ITE on the training and test set in terms of the risk difference (see Equation 2.11). For visual validation, we show the densities of the estimated ITEs on both datasets, and the ITE-ATE plots to assess whether the estimated ITEs align with the observed outcomes.

3.2.3 Results

In this section, we present the results of the ITE estimation on the International Stroke Trial (IST) dataset.

The observed average treatment effect (ATE), defined as $P(Y = 1|T = 1) - P(Y = 1|T = 0)$, was -2.4% absolute risk reduction on the training set, with a 95% confidence interval from -4.1% to -0.6%. The interval was computed using the Wald method for risk differences, with standard error $\sqrt{p_1(1-p_1)/n_1 + p_0(1-p_0)/n_0}$, where p_1 and p_0 are the event rates in the treated and control groups. On the test set, the observed treatment effect was -0.1%, with a 95% confidence interval from -2.6% to 2.3%.

To estimate ITEs (Equation 2.11), we applied three models: a T-learner logistic regression, a T-learner tuned random forest, and an S-learner TRAM-DAG. The results for each model are shown in terms of (1) the density of predicted ITEs, and (2) ITE-ATE plots, which display the empirical risk difference within subgroups of estimated ITEs, including 95% confidence intervals. Results are presented in Figures 3.9 (logistic regression), 3.10 (random forest), and 3.11 (TRAM-DAG). Additionally, calibration plots are provided in Appendix 5.7, Figures 5.5 - 5.7.

The estimated average treatment effect on the test set, calculated as $ATE_{\text{pred}} = \text{mean}(ITE_{\text{pred}})$, was -2.5% for the T-learner logistic regression, -2.2% for the T-learner tuned random forest, and -3.1% for the S-learner TRAM-DAG.

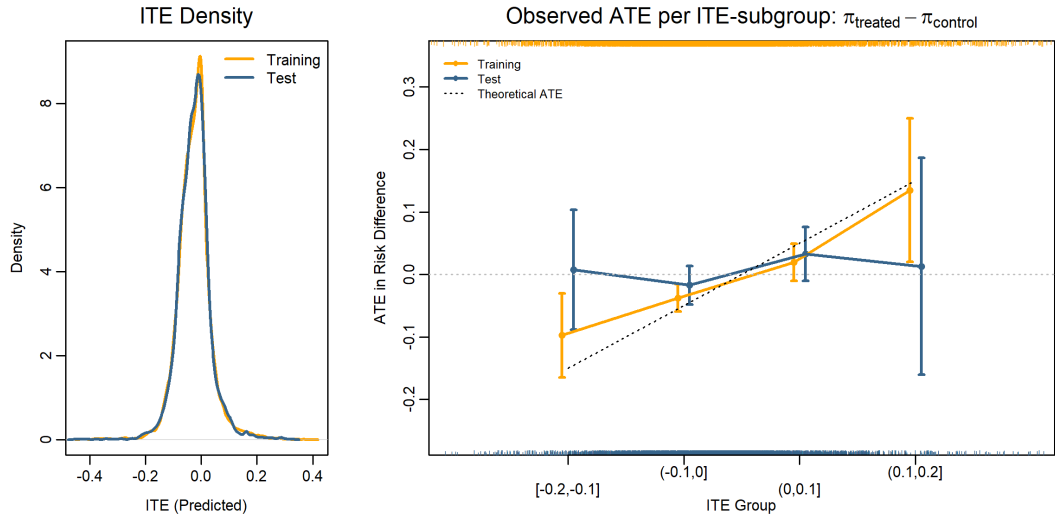


Figure 3.9: Results for the International Stroke Trial (IST) using the T-learner logistic regression. Left: density of predicted ITEs in the training and test sets; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

enforce that starts after all floats have been displayed

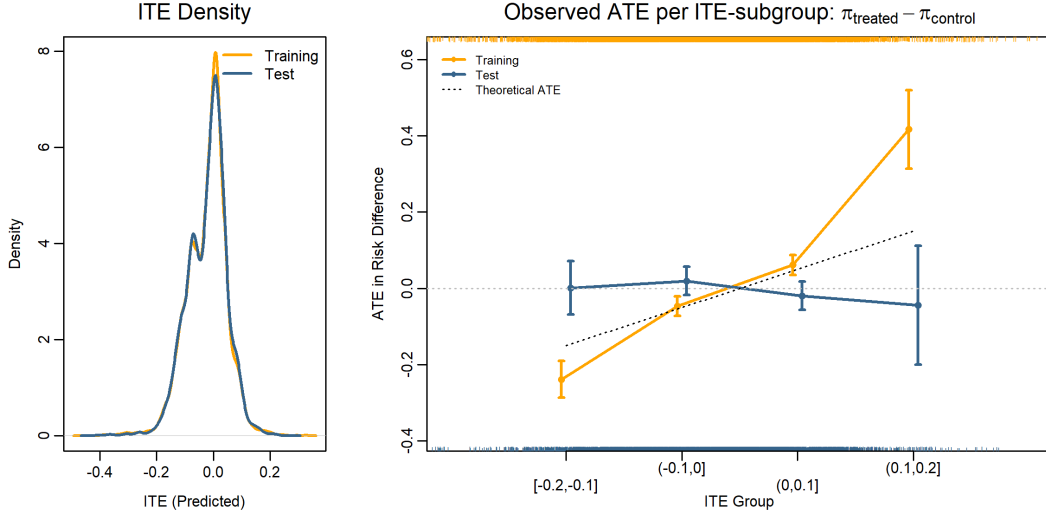


Figure 3.10: Results for the International Stroke Trial (IST) using the T-learner tuned random forest. Left: density of predicted ITEs in the training and test sets; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

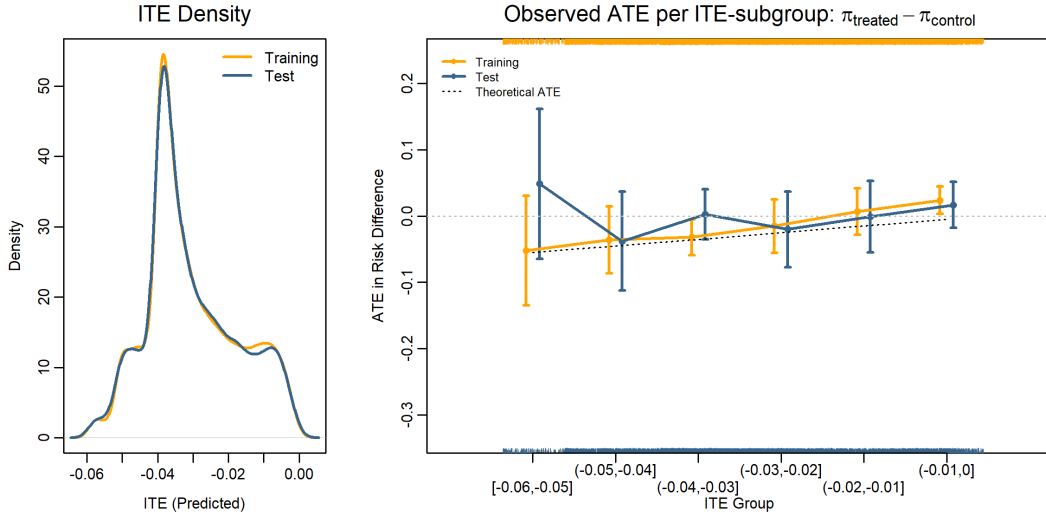


Figure 3.11: Results for the International Stroke Trial (IST) using the S-learner TRAM-DAG. Left: density of predicted ITEs in the training and test sets; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

3.2.4 Discussion of Experiment 2

We observed similar results to those reported by [Chen *et al.* \(2025\)](#) when estimating ITEs on the International Stroke Trial dataset across all three models: the T-learner logistic regression, the T-learner tuned random forest, and the S-learner TRAM-DAG. The logistic model showed moderate discrimination in the training set, which did not generalize to the test set, as illustrated by the ITE-ATE plot in Figure 3.9. The tuned random forest model showed stronger discrimination in the training set but similarly failed to generalize to the test set (see Figure 3.10). In contrast, the S-learner TRAM-DAG estimated less heterogeneity than the other two models, as shown in the density plot in Figure 3.11, resulting in weak discrimination in both the training and test sets, visible in the corresponding ITE-ATE plot. For all three models, the confidence intervals in the ITE-ATE plots on the test set included the zero line, suggesting no significant effect in

any of the estimated ITE subgroups.

Poor calibration does not appear to explain the limited ITE performance, as calibration on the test set was good, as shown in Appendix 5.7, Figures 5.5-5.7. However, since the ground truth is unknown, it remains unclear whether the models fail to detect true treatment effect heterogeneity, or whether the heterogeneity is too small, or driven by unobserved effect modifiers. We investigate these possibilities further in Experiment 3, Section 3.3.

3.3 Experiment 3: ITE model robustness in RCTs (simulation)

3.3.1 Motivation

In this section, we perform a simulation study to estimate ITEs using different models in an RCT setting under various scenarios. The aim is to identify conditions under which ITE estimation fails, and whether such failure is model-agnostic – i.e., driven by external factors such as unobserved covariates or the strength of the treatment effect, rather than by the model class itself. This may provide insight into why ITE estimation can fail in real-world applications, as demonstrated by [Chen *et al.* \(2025\)](#) on the IST dataset and replicated in our own work in Experiment 2 (Section 3.2).

3.3.2 Setup

The simulation is based on a data-generating process (DGP) that resembles an RCT. We assume a binary outcome and a set of covariates that influence the outcome. There may also be treatment-covariate interactions that are responsible for heterogeneity in the treatment effect.

Data-generating process: Data is generated similarly to the approach proposed by [Hoogland *et al.* \(2021\)](#). The binary treatment (T) is sampled from a Bernoulli distribution with probability 0.5. The five covariates (\mathbf{X}), representing patient-specific characteristics at baseline, are drawn from a multivariate standard normal distribution with a compound symmetric covariance matrix ($\rho = 0.1$). The binary outcome (Y) is sampled from a Bernoulli distribution with probability $P(Y_i = 1 \mid \mathbf{X} = \mathbf{x}_i, T = t_i) = \text{logit}^{-1}(\beta_0 + \beta_T t_i + \beta_X^\top \mathbf{x}_i + t_i \cdot \beta_{TX}^\top \mathbf{x}_{TX,i})$, where i denotes the patient index, and $\mathbf{x}_{TX,i}$ denotes the subset of covariates that interact with the treatment.

The simulated datasets are generated under three scenarios, where coefficients are set to different values or not all variables are observed. In Scenario 3.1, the coefficients are: $\beta_0 = 0.45$ (intercept), $\beta_T = -0.85$ (direct treatment effect), $\beta_X = (-0.5, 0.8, 0.2, 0.6, -0.4)$ (direct covariate effects), and $\beta_{TX} = (0.9, 0.1)$ (interaction effects between treatment and covariates X_1 and X_2 on the outcome). In Scenario 3.2, the same coefficients are used, but the covariate X_1 , which is responsible for a large portion of the heterogeneity, is not observed in the final dataset. This is expected to cause difficulties in estimating the ITE. In Scenario 3.3, the coefficients for the direct treatment and interaction effects are set to $\beta_T = -0.05$ and $\beta_{TX} = (-0.01, 0.03)$ to represent a weak treatment effect and low heterogeneity. All other coefficients remain unchanged, and all covariates are observed. The DAGs corresponding to the three scenarios are presented in Figure 3.12.

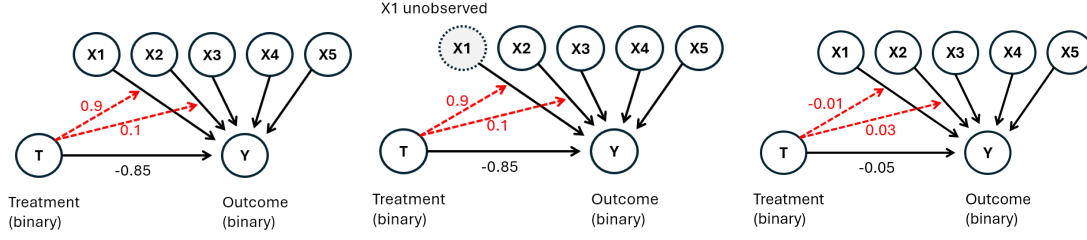


Figure 3.12: Data-generating process (DGP) for the three scenarios in the ITE simulation study (RCT). Interaction effects between treatment (T) and covariates (X_1 and X_2) on the outcome (Y) are shown in red. Left: Scenario 3.1, where all covariates are observed and both treatment effect and heterogeneity are strong; Middle: Scenario 3.2, with the same DGP as in Scenario 3.1, but where covariate X_1 is not observed; Right: Scenario 3.3, where the treatment effect and heterogeneity are weak, and all covariates are observed.

Models for ITE estimation: We applied the following models in R to estimate individualized treatment effects (ITEs) according to Equation 2.11: T-learner logistic regression (`stats` package), T-learner and S-learner logistic lasso regression (`glmnet` package (Friedman *et al.*, 2010), with λ selected via 10-fold cross-validation), T-learner random forest (`randomForest` package (Breiman, 2001), 100 trees), and T-learner tuned random forest (`comets` package (Kook, 2024), which tunes `mtry` and `max.depth` using out-of-bag error, 500 trees).

While all models were run, we focus on two for the main results: the T-learner logistic regression as a baseline (matching the DGP), and the T-learner tuned random forest as a more flexible non-parametric model. Results for the default random forest in Scenario 3.1 are shown in Appendix 5.8 to highlight the role of model tuning in avoiding overfitting and ensuring proper calibration.

Each model was trained and evaluated on independent datasets of 10,000 samples generated from the same DGP. Although TRAM-DAGs are well suited for ITE estimation, we did not include them here, as the goal of this experiment is to compare simpler and more complex models under varying conditions. Furthermore, when specified as T-learner with standard logistic latent distribution and linear shift effects of covariates, the TRAM-DAG would be equivalent to a logistic regression for the binary outcome.

Model evaluation: Model performance is evaluated visually on both the training and test datasets. For predictive performance, we present true vs. predicted probabilities $P(Y = 1 | X, T)$ to assess how well each model is calibrated. Plots of true vs. predicted ITEs (calculated according to Equation 2.11) show how closely the model estimates match the true effects. Since the true probabilities and ITEs are known by design in this simulation, direct evaluation of calibration and prediction accuracy is possible, unlike in real-world applications.

To assess whether estimated ITEs correspond to actual observed outcomes, we use ITE-ATE plots. These show the observed average treatment effect (ATE), calculated as $P(Y = 1 | T = 1) - P(Y = 1 | T = 0)$, in the respective subgroups of estimated ITEs. Accurate models should produce ITE-ATE points that align with the identity line.

These simulation scenarios allow us to assess ITE estimation performance under challenging conditions such as omitted variables and weak treatment effects. The subsequent results reveal which models remain robust under such violations and provide insight into possible real-world estimation failures.

3.3.3 Results and discussion

In this section, we present the performance of two causal machine learning models – T-learner logistic regression and T-learner tuned random forest – for estimating ITEs under the three

simulated scenarios introduced in Section 3.3.2. Scenario 3.1 (Section 3.3.3.1) represents the ideal case where all covariates are observed and both treatment and interaction effects are strong. In Scenario 3.2 (Section 3.3.3.3), a key effect modifier is unobserved, and in Scenario 3.3 (Section 3.3.3.3), treatment and interaction effects are weak, but all covariates are observed.

3.3.3.1 Scenario 3.1: Fully observed, large effects

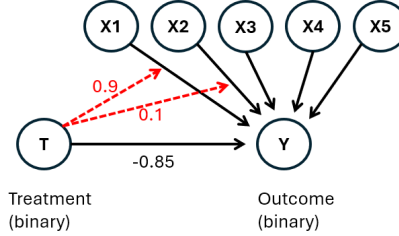


Figure 3.13: DAG for Scenario 3.1, where all variables are observed and both treatment and interaction effects are strong. This DAG was previously shown in Figure 3.12 and is re-plotted here for convenience. The numbers indicate the coefficients on the log-odds scale. Red arrows represent interaction effects between treatment (T) and covariates (X_1 and X_2) on the outcome (Y).

In Scenario 3.1 (Figure 3.13), where treatment effect heterogeneity was large and all covariates were observed, the T-learner logistic regression accurately estimated the ITE. The observed ATE, conditional on the respective ITE subgroup, was well calibrated in both the training and test datasets (Figure 3.14). This is as expected, since the data were generated with the same model class (logistic regression), and applying logistic regression as a T-learner for ITE estimation can accurately capture the interaction effects.

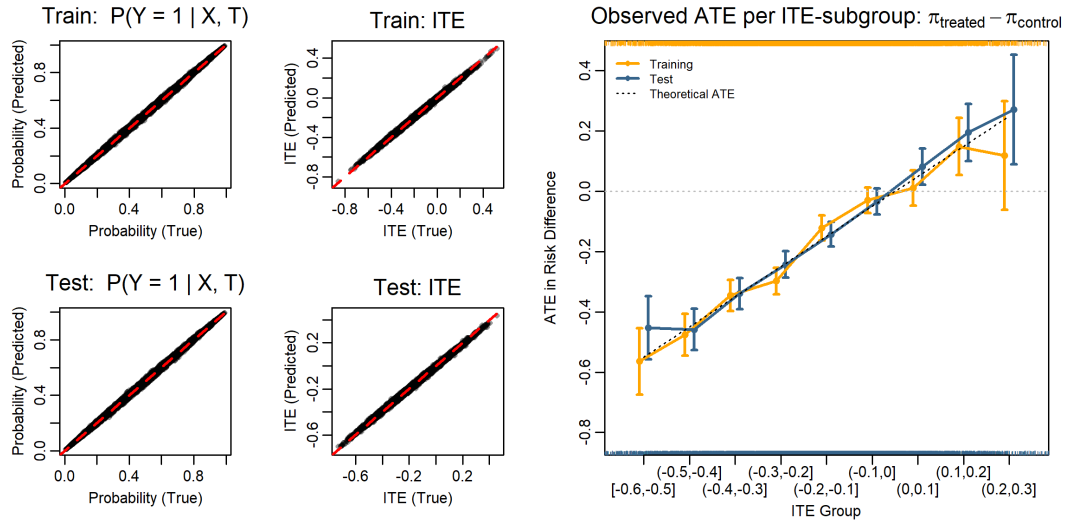


Figure 3.14: Results of the T-learner logistic regression in Scenario 3.1, where the DAG is fully observed and both treatment and interaction effects are strong. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

The tuned random forest model also performed well (see Figure 3.15), but less accurate compared to the logistic model (see Figure 3.14). Choosing a different model class than that used in the DGP may lead to worse prediction accuracy in terms of $P(Y = 1 | X, T)$ and ITE.

This difference between the two models arises because the logistic regression model has only a small number of parameters, and with sufficient data, these parameters can converge to their true values as used in the logistic DGP, allowing near-perfect recovery of the true probabilities and thus ITEs. In contrast, the non-parametric random forest must infer the underlying probabilities from the observed binary outcomes (0 or 1), which are themselves realizations of a Bernoulli process. This introduces inherent noise, making it harder for the model to estimate the true risk accurately – even with large sample sizes. Nonetheless, the tuned random forest also captured the general trend of the ITEs, as reflected in the ITE-ATE plot, Figure 3.15. Both models were able to capture treatment effect heterogeneity well under full observability of covariates.

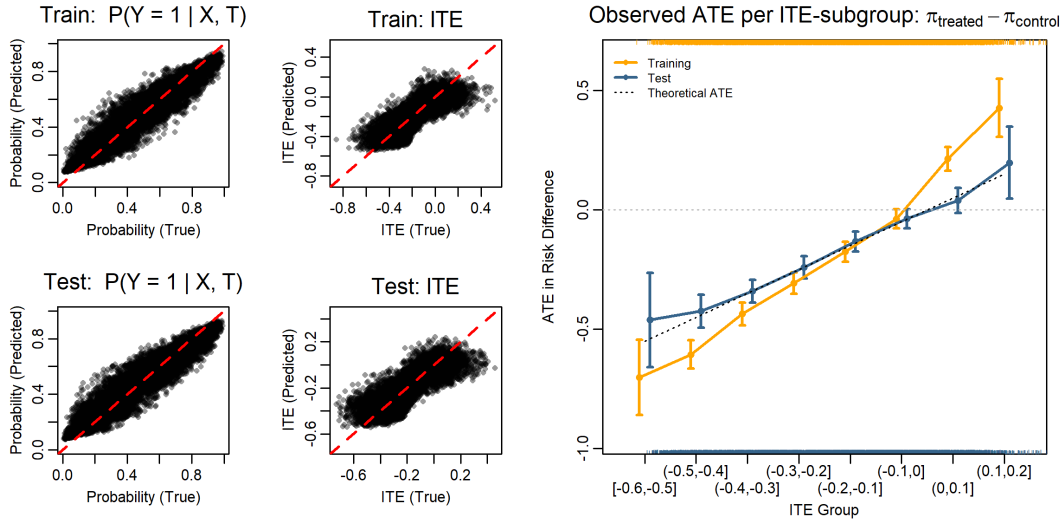


Figure 3.15: Results of the T-learner tuned random forest in Scenario 3.1, where the DAG is fully observed and both treatment and interaction effects are strong. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

In contrast, the default random forest (i.e., without hyperparameter tuning) performed worse than its tuned counterpart (see Appendix 5.8). As shown in the corresponding ITE-ATE plot, Figure 5.9, the model exhibited poor calibration and inaccurate ITE estimates, highlighting the importance of proper tuning to ensure reliable ITE estimation and avoid overfitting.

3.3.3.2 Scenario 3.2: Unobserved interaction

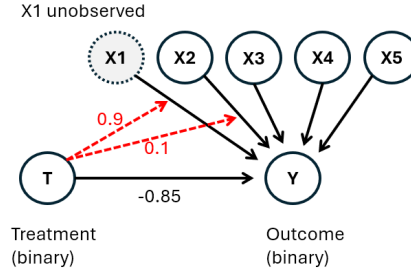


Figure 3.16: DAG for Scenario 3.2, where there are strong treatment and interaction effects, but variable X_1 is not observed. This DAG was previously shown in Figure 3.12 and is re-plotted here for convenience. The numbers indicate the coefficients on the log-odds scale. Red arrows represent interaction effects between treatment (T) and covariates (X_1 and X_2) on the outcome (Y).

In Scenario 3.2 (Figure 3.16), where treatment effect heterogeneity remained large but one important interaction covariate (X_1) was not observed, prediction accuracy decreased for both models, and the estimated heterogeneity in terms of the ITE was smaller than the true heterogeneity (see true vs. estimated ITE plots in Figures 3.17 and 3.18).

Although a considerable number of individuals had a true ITE that was positive (visible in the true vs. estimated ITE plot in Figure 3.17), the T-learner logistic regression did not predict a single positive ITE. This shows that the missing covariate X_1 in this example prevents detection of individuals who would actually benefit from the treatment. In practice, where the true ITEs are unknown, we would evaluate the ITE estimation in terms of the ITE-ATE plot (Figure 3.17), which in this case would suggest good calibration of ITEs, as they follow the identity-line. However, we still did not identify the patients that would actually benefit from the treatment. This implies that even if the ITE-ATE plot suggests heterogeneity and that ITE estimates are well calibrated, there may still be additional heterogeneity, which was not captured by the model.

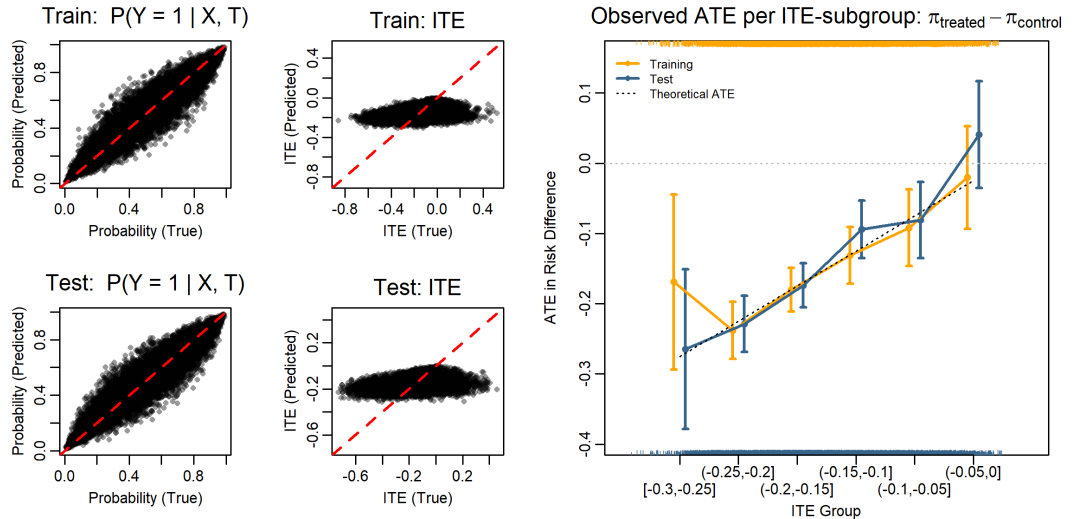


Figure 3.17: Results of the T-learner logistic regression in Scenario 3.2, where there are strong treatment and interaction effects, but variable X_1 is not observed. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

In contrast, the T-learner tuned random forest (Figure 3.18) estimated larger treatment effect heterogeneity than the logistic model, but still could not accurately estimate the ITE and also failed to detect patients who would benefit from the treatment. The ITE-ATE plot in Figure 3.18 illustrates that the model discriminates too strongly in the training set and does not generalize well to the test set.

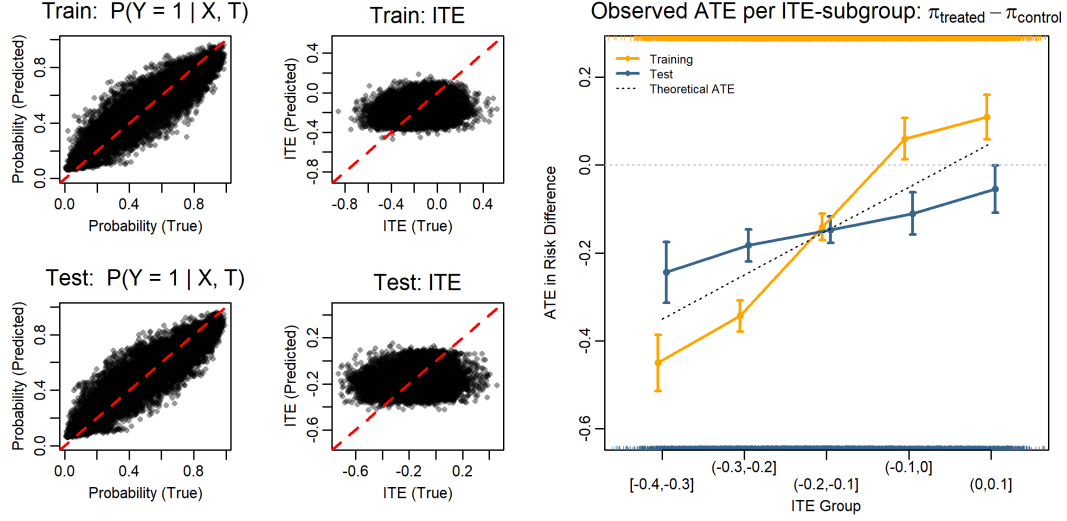


Figure 3.18: Results of the T-learner tuned random forest in Scenario 3.2, where there are strong treatment and interaction effects, but variable X_1 is not observed. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

3.3.3.3 Scenario 3.3: Fully observed, small effects

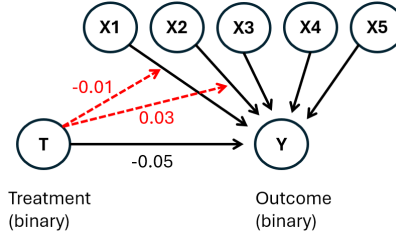


Figure 3.19: DAG for Scenario 3.3, where all variables are observed and both treatment and interaction effects are weak. This DAG was previously shown in Figure 3.12 and is re-plotted here for convenience. The numbers indicate the coefficients on the log-odds scale. Red arrows represent interaction effects between treatment (T) and covariates (X_1 and X_2) on the outcome (Y).

In Scenario 3.3 (Figure 3.19), where the true treatment effect heterogeneity was small and all covariates were observed, the T-learner logistic regression estimated a larger heterogeneity than actually present (see true vs. estimated ITE plots in Figure 3.20). In the ITE-ATE plot in Figure 3.20, the confidence intervals of all ITE subgroups overlap and include the zero line, indicating that the treatment effect is not significantly different from zero. This matches expectations given the small true effect sizes.

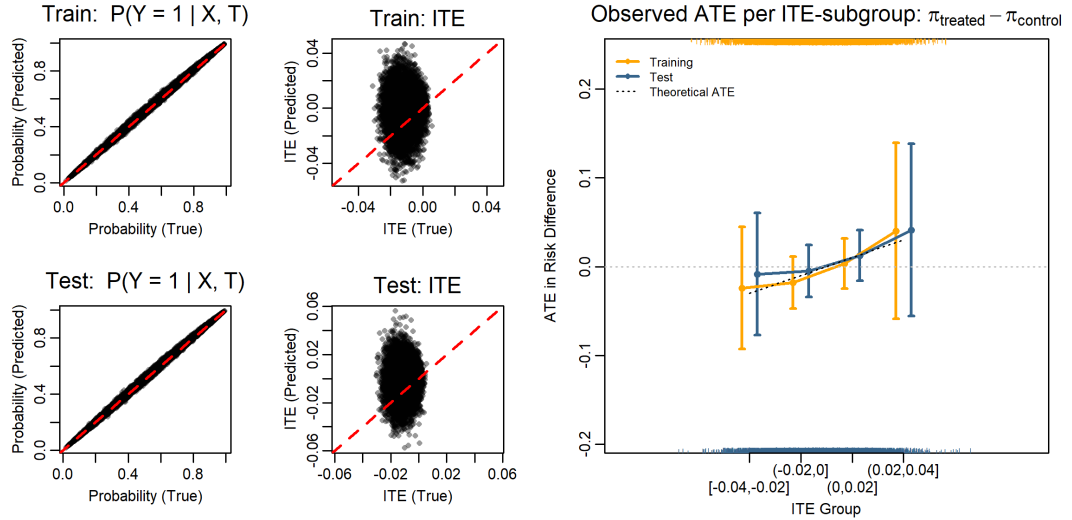


Figure 3.20: Results of the T-learner logistic regression in Scenario 3.3, where the DAG is fully observed and both treatment and interaction effects are weak. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

However, the T-learner tuned random forest model incorrectly estimated even larger treatment effect heterogeneity than the logistic regression model (see true vs. estimated ITE plots in Figure 3.21). As shown in the ITE-ATE plot in Figure 3.21, the model exhibited strong discrimination in the training set but did not replicate this pattern in the test set, where – regardless of the estimated ITE – the observed outcomes were similar across all subgroups. These results indicate that when evaluating the ITE predictions on the test set, it correctly suggests that no significant treatment effect heterogeneity is present.

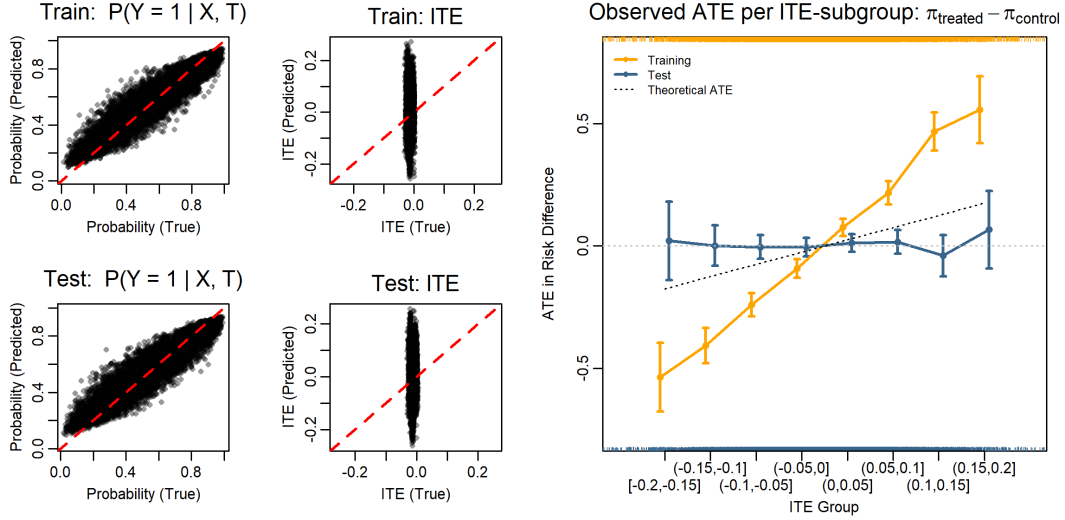


Figure 3.21: Results of the T-learner tuned random forest in Scenario 3.3, where the DAG is fully observed and both treatment and interaction effects are weak. Left: true vs. predicted probabilities for $P(Y = 1 | X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

3.3.4 Discussion of Experiment 3

Tuning more flexible models like random forests using cross-validation improved generalization to the test set but led to poor calibration in terms of predicted probabilities vs. empirically observed outcomes in the training set. An illustrative case is shown in Appendix 5.9 for the T-learner tuned random forest in Scenario 3 (with weak effects; Section 3.3.3.3), where calibration was poor in the training set but aligned well with the identity line in the test set. We repeatedly observed this pattern in the tuned random forest when, in the ITE-ATE plot, results from the training set did not generalize to the test set (e.g. in Figures 3.18 and 3.21). This highlights the importance of evaluating models on an independent test set, when tuning a model to prevent overfitting, although evaluation on a test set should be done in any case.

In this experiment, we showed that even when causal ML models for ITE estimation are well calibrated in terms of prediction accuracy $P(Y = 1 | \mathbf{X}, T)$, they can still fail to estimate the ITE accurately under less favorable scenarios. In cases of full observability of covariates but low interaction effects (Scenario 3.3; Section 3.3.3.3), models may estimate too high heterogeneity that is not present in the data. However, this can become visible in the ITE-ATE plot on the test set, which reveals that the apparent heterogeneity does not generalize. But we also observed that when important effect-modifying covariates are missing (Scenario 3.2; Section 3.3.3.2), the models may fail to detect treatment effect heterogeneity altogether. In such cases, the estimated ITEs may be too small or even negative, suggesting that the model does not capture the true treatment effect heterogeneity. This makes it difficult to distinguish between a true lack of heterogeneity and the failure to capture it due to unobserved effect modifiers.

[Vegetabile \(2021\)](#) also analyzed the effect of unobserved interaction variables. He pointed out that as long as all confounding variables \mathbf{X} are observed and conditioned on, the ignorability assumption (Equation 2.8) required for ITE estimation is satisfied – even in the presence of an unobserved interaction variable Z . However, if such a variable Z exists, the estimated ITEs would be biased, and this issue could arise even in an RCT setting where confounding is removed through randomization. This aligns with our observations in Scenario 3.2 (Section 3.3.3.2).

[Nichols \(2007\)](#) discusses various methods for estimating causal effects from observational data, including in the presence of unobserved variables. One of these methods, instrumental

variables (IV), can help reduce bias from unobserved confounding. Whether IV methods can also address unobserved effect modifiers in the context of ITE estimation is not something we explored, and remains beyond the scope of this thesis.

3.4 Experiment 4: ITE estimation with TRAM-DAGs (simulation)

3.4.1 Motivation

We claim that the TRAM-DAG framework can be effectively used for unbiased ITE estimation on observational data that includes both confounding and mediating variables – provided that identifiability assumptions are met and the underlying DAG is fully known. To evaluate this, we apply TRAM-DAGs in both confounded and randomized settings, using data simulated according to the DAGs shown in Figure 3.22. The binary treatment variable (X_4) represents the intervention, and the goal is to estimate the ITE for a continuous outcome Y .

To assess model performance under varying conditions, we conduct this experiment across three scenarios:

- **Scenario 4.1** (Section 3.4.3.1): both direct and interaction effects of treatment are present,
- **Scenario 4.2** (Section 3.4.3.2): only a direct treatment effect is included,
- **Scenario 4.3** (Section 3.4.3.3): only interaction effects are included.

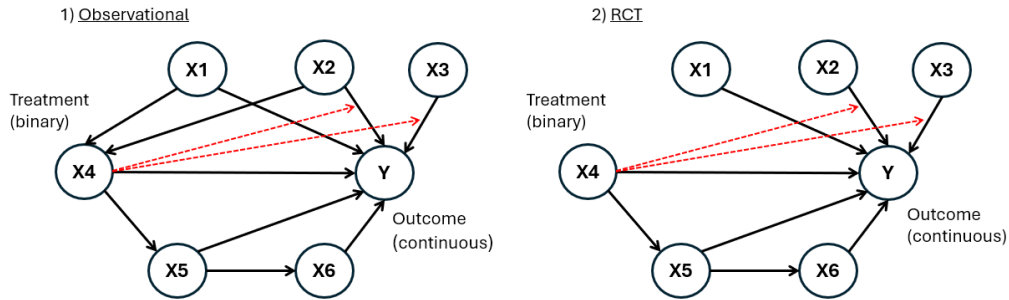


Figure 3.22: DAGs used for the simulation to estimate the ITE. Left: Observational; Right: RCT setting. The source nodes X_1 , X_2 , and X_3 come from a multivariate standard normal distribution ($\rho = 0.1$). In the observational setting, the binary treatment X_4 depends on the parents X_1 and X_2 . In the RCT setting, this dependency is omitted due to randomization. The outcome Y depends on all variables, with additional interaction effects between the treatment and the variables X_2 and X_3 . All variables except the treatment X_4 are continuous.

Illustrative scenario: A possible real-world scenario that follows the structure of the proposed DAG could be the following: A marketing campaign is conducted to increase customer spending. The treatment is the marketing email (X_4) sent to customers. If the treatment is not randomized, it depends on prior total spend (X_1) and the customer engagement score (X_2). The outcome is the total spend in the 30 days following the email, denoted as Y . The past total spend (X_1) and customer engagement score (X_2) act as confounders, influencing both the treatment and the outcome. The customer satisfaction score (X_3), obtained from a recent survey, is another predictor. The time spent on the website after receiving the email (X_5) is a mediator that affects the number of product pages viewed (X_6), which in turn influences the

total spend (Y). Interaction effects exist between the treatment (X_4) and both X_2 and X_3 , meaning the treatment effect differs based on the customer’s engagement and satisfaction levels. The goal is to estimate the individualized treatment effect (ITE) of the marketing email (X_4) on the total spend (Y), in order to personalize customer targeting.

3.4.2 Setup

Data-generating process: The standard logistic distribution was chosen as the noise distribution to align with other examples in this thesis. Any other noise distribution could also be used here, as we are not interested in coefficient interpretability in this experiment. All variables except the binary treatment X_4 are continuous. The source nodes X_1 , X_2 , and X_3 are generated from a multivariate standard normal distribution with a compound symmetric covariance matrix ($\rho = 0.1$). These variables represent baseline patient characteristics.

In the observational setting, X_1 and X_2 act as confounders by influencing both the treatment assignment X_4 and the outcome Y . In the RCT setting, these dependencies are removed due to randomization. The mediator X_5 depends on treatment X_4 , and X_6 depends on X_5 . The log-odds of the continuous outcome Y depend linearly on all covariates, including additional interaction terms between the treatment and X_2 and X_3 . Equation 3.1 defines the outcome Y on the log-odds scale:

$$h(y \mid \mathbf{X}) = h_I(y) + \beta_X^\top \mathbf{X} + X_4 \cdot (\beta_{TX}^\top \mathbf{X}_{TX}) \quad (3.1)$$

Here, $h_I(y)$ is the intercept function, \mathbf{X} is the full covariate vector, and $\mathbf{X}_{TX} = \{X_2, X_3\}$ denotes the interaction covariates that affect the outcome only when treatment is applied ($X_4 = 1$). The intercept function $h_I(y)$ must be smooth and monotonically increasing. We define it as $h_I(y) = \tan(y/2)/0.2$ for $y \in [-2, 2]$, and extrapolate linearly at the boundaries.

The coefficients are set as $\beta_X = (-0.5, 0.5, 0.2, 1.5, -0.6, 0.4)$, where the value 1.5 represents the direct effect of treatment X_4 on the outcome. The interaction coefficients are set to $\beta_{TX} = (-0.9, 0.7)$.

Three scenarios: The experiment is conducted under three different scenarios regarding the effect of the treatment on the outcome Y in the DGP: Scenario 4.1 (Section 3.4.3.1) with both direct treatment and interaction effects, Scenario 4.2 (Section 3.4.3.2) with only a direct effect, and Scenario 4.3 (Section 3.4.3.3) with only interaction effects. Depending on the scenario, the corresponding coefficients in β_X and β_{TX} in Equation 3.1 are set to zero.

TRAM-DAG estimation: In both the observational and RCT settings, the TRAM-DAG is fitted as an S-learner (i.e., a single model including the treatment variable). To allow for full flexibility, all nodes with parents are modeled using complex intercepts with three hidden layers of size (10, 10, 10), without applying batch normalization or dropout, using ReLU activation. This architecture enables the model to learn nonlinearities and interactions between the treatment and covariates, allowing it to estimate both potential outcomes. The model is trained on a dataset of 20,000 samples. To prevent overfitting, an additional validation set of 10,000 samples is used, and the final model is selected using early stopping based on the validation loss.

ITE estimation procedure: In contrast to much of the literature we reviewed, where ITEs are typically defined in terms of expected values of potential outcomes (Equation 2.10), we estimate the quantile treatment effect (QTE), specifically at the median (Equation 3.2). For each individual, we calculate the difference between the medians of the potential outcome distributions under treatment and control. Note that estimating the potential outcomes in terms of expected values would also be possible – either by repeatedly sampling from each outcome distribution or, potentially, by numerical integration. However, for this experiment, we chose to estimate the QTE. Chernozhukov and Hansen (2005), for example, highlighted the ability of quantile regression models in heterogeneous treatment effect estimation. QTEs are particularly

relevant when the distributional behavior of outcomes beyond the mean is of interest. The median QTE is defined as

$$\text{QTE}^{(0.5)}(\mathbf{x}) = Q_{Y(1)|\mathbf{x}=\mathbf{x}}(0.5) - Q_{Y(0)|\mathbf{x}=\mathbf{x}}(0.5), \quad (3.2)$$

where $Q_{Y(t)|\mathbf{x}=\mathbf{x}}(q)$ denotes the q -th quantile of the potential outcome distribution under treatment t .

Once the TRAM-DAG model is fitted on observed data, we can access the inverse transformation functions $X_i = h^{-1}(Z_i | \text{pa}(X_i))$, which represent the structural equations of the DAG. ITE estimation proceeds in three steps, according to Algorithm 3. First, the latent values z_{ij} for the explanatory variables $X_i \in \{X_1, X_2, X_3, X_5, X_6\}$ are computed for each sample j using the transformation functions conditioned on their observed parents. Second, the treatment variable X_4 is intervened on using the do-operator for both $X_4 = 0$ and $X_4 = 1$. For each of these two treatment states, X_5 , X_6 , and the potential outcome (distribution) Y are sampled sequentially using the latent encodings and inverse transformations. This means that the counterfactuals for X_5 and X_6 are determined. This results in two potential outcome distributions per individual, as illustrated in Figure 3.23. Finally, for each individual, the median is determined for both potential outcome distributions and the QTE is calculated as the difference between the two medians (Equation 3.2). For simplicity, we will refer to QTEs as ITEs throughout the remainder of the experiment.

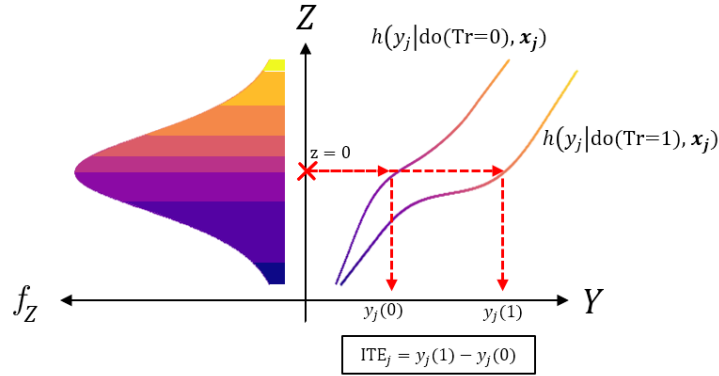


Figure 3.23: ITE estimation in terms of quantile treatment effect (QTE) (Equation 3.2) at the median with TRAM-DAGs. The two transformation functions represent the distributions of the potential outcomes under both treatments. For the $\text{QTE}(0.5)$, the median of the latent distribution (0 for the standard logistic) is evaluated on both transformation functions to determine the median potential outcomes. We define the ITE for an individual as the difference of the median potential outcomes. Plot adapted from visualizations similar to those used in Sick and Dürr (2025).

Algorithm 3 ITE estimation (QTE) using TRAM-DAGs

Input: Fitted TRAM-DAG, dataset of n individuals

for each individual $j = 1$ to n **do**

Step 1: Determine latent values

for each explanatory node $X_i \in \{X_1, X_2, X_3, X_5, X_6\}$ **do**

 Compute latent value: $z_{ij} = h_i(x_{ij} \mid \text{pa}(x_{ij}))$

end for

Step 2: Generate potential outcomes under treatment and control

for $x_4 \in \{0, 1\}$ **do** ▷ Simulate both treatment states

 Fix $X_4 = x_4$ (intervention)

 Sample X_5 and X_6 sequentially using z_{ij} and inverse transformations

 Sample potential outcome $y_j(x_4)$ using $z_{7,j} = 0$ (median of the potential outcome distribution)

end for

Step 3: Compute ITE (QTE) for individual j

$\text{ITE}_j = y_j(1) - y_j(0)$

end for

Output: ITE estimates $\{\text{ITE}_j\}_{j=1}^n$

Model evaluation: Validation is conducted on the training dataset and on an independent test dataset, each of which consists of 20,000 samples. During the data-generating process, the true potential outcomes under both treatment states were recorded for each individual, which allows for exact computation of the true ITE. The estimated ITEs are evaluated against the true values using several visual and numerical metrics. These include density plots of the estimated ITEs (e.g. Figure 3.28), scatter plots of true vs. estimated ITEs (e.g. Figure 3.29), and ITE-ATE plots (e.g. Figure 3.30) where the observed ATE per ITE subgroup is computed as the difference in medians. In addition, the average of the estimated ITEs is compared to the true average ITE and to the empirical ATE from the RCT setting (e.g. Table 3.1). Among the evaluated metrics, the scatter plot of true versus estimated ITEs provides the most direct insight into estimation accuracy.

3.4.3 Results and Discussion

First, we present the results for Scenario 4.1 (Section 3.4.3.1), which includes both direct and interaction effects of treatment. Then, we show the results for Scenario 4.2 (Section 3.4.3.2), which has a direct effect but no interaction effects, and finally Scenario 4.3 (Section 3.4.3.3), which includes interaction effects but no direct treatment effect. For each scenario, we compare the results in an observational setting with confounded treatment allocation and in a randomized controlled trial (RCT) setting without confounding. We also compare the average treatment effect (ATE), which can be directly calculated in the RCT setting on observed outcomes, with the ATE derived from the estimated ITEs. All ITEs presented in this experiment are technically quantile treatment effects (QTEs), as defined in Equation 3.2, based on the medians of the potential outcomes. For simplicity, we will refer to them as ITEs. The aim is to investigate how the TRAM-DAG performs in the presence or absence of direct and interaction effects of the treatment, both in confounded and in randomized settings, for the purpose of ITE estimation.

3.4.3.1 Scenario 4.1: Direct and interaction effects

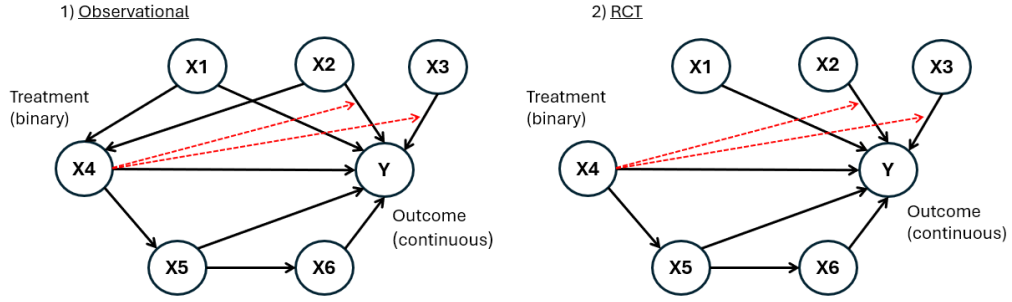


Figure 3.24: DAGs for Scenario 4.1, which includes a direct effect of the treatment on the outcome and additional interaction effects with covariates X_2 and X_3 . These DAGs were previously shown in Figure 3.22 and are re-plotted here for convenience. Left: Observational; Right: RCT setting.

Scenario 4.1 includes a direct effect of the treatment on the outcome, and interaction effects with X_2 and X_3 , as shown in Figure 3.24.

In the observational setting, treatment was confounded by X_1 and X_2 . In the train set, 38.6% of individuals were in the control group and 61.4% in the treatment group. The test set had a similar distribution.

In the RCT setting, treatment was randomly assigned. In the train set, 49.8% were in control and 50.2% in treatment; in the test set, the shares were 50.2% and 49.8%, respectively.

Figure 3.25 shows the true ITE distribution resulting from the DGP, which displays some heterogeneity due to interaction effects. Figure 3.26 shows the marginal distributions of all variables in the DGP and as estimated by the TRAM-DAG. The distribution of the outcome Y under $\text{do}(X_4 = 0)$ and $\text{do}(X_4 = 1)$ is shown in Figure 3.27. The ITEs were estimated as the difference in medians of the potential outcomes. Figure 3.28 compares the densities of the estimated and true ITEs. Figure 3.29 shows the scatterplots of estimated vs. true ITEs. In both observational and RCT settings, the estimated ITEs are close to the true ones in both train and test sets. Figures 3.30 and 3.31 show the ITE-ATE plots, where the ATE is calculated as the median difference in observed outcomes within ITE subgroups. The trends are similar across train and test sets and follow the calibration line.

The ATEs calculated based on different measures in both the training and test sets are shown in Table 3.1. In the RCT setting (training set), the difference in means of the outcomes between the two treatment groups was -0.563 , with a 95% Wald confidence interval from -0.582 to -0.543 . Note that the ATE in terms of difference in means cannot be directly compared to the ATE based on difference in medians.

Table 3.1: Scenario 4.1, including direct and interaction effects: Comparison of ATE measures across training and test sets for the observational and RCT settings. $Y_{\text{observed}}^{(\text{Tr})}$ denotes the observed outcome under the treatment (Tr) actually received. Estimates based on these observed outcomes (means and medians) are provided only for the RCT setting, as the observational setting is confounded. The true ITEs (ITE_{true}) were calculated for each individual based on the data-generating process. In contrast, the estimated ITEs ($\text{ITE}_{\text{estimated}}$) were obtained from the TRAM-DAG trained on observed data. The estimated ATE from $\text{mean}(\text{ITE}_{\text{estimated}})$ can be directly compared to the true $\text{mean}(\text{ITE}_{\text{true}})$, whereas comparisons to empirical ATEs from observed outcome differences should be interpreted with caution. All ITEs were computed as quantile treatment effects (QTEs) based on the median of the potential outcome distributions, as defined in Equation 3.2.

Measure	Observational		RCT	
	Train	Test	Train	Test
ATE as $\text{mean}(Y_{\text{observed}}^{(1)}) - \text{mean}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.563	-0.563
ATE as $\text{median}(Y_{\text{observed}}^{(1)}) - \text{median}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.626	-0.638
ATE as $\text{mean}(\text{ITE}_{\text{true}})$	-0.620	-0.622	-0.620	-0.622
ATE as $\text{mean}(\text{ITE}_{\text{estimated}})$	-0.617	-0.620	-0.619	-0.622

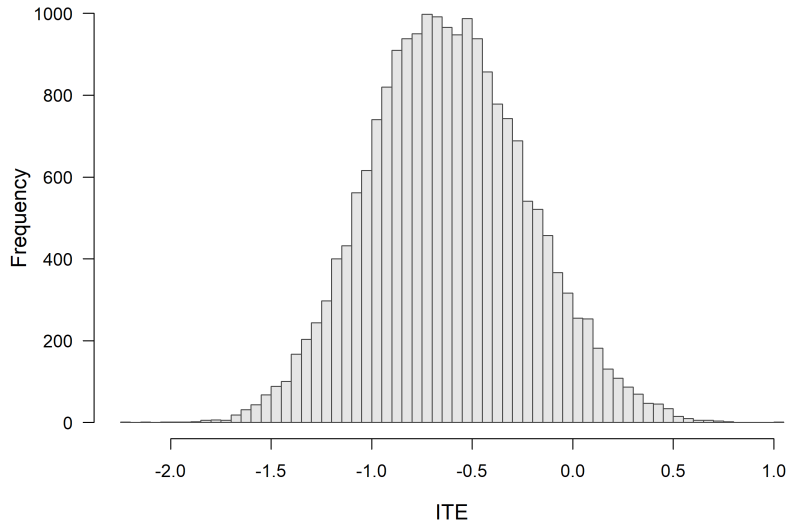


Figure 3.25: True ITE distribution resulting from the DGP for Scenario 4.1, which includes both direct and interaction effects. The true ITEs are identical for each individual in the observational and RCT settings, as they are based on the potential outcomes under both treatment conditions.

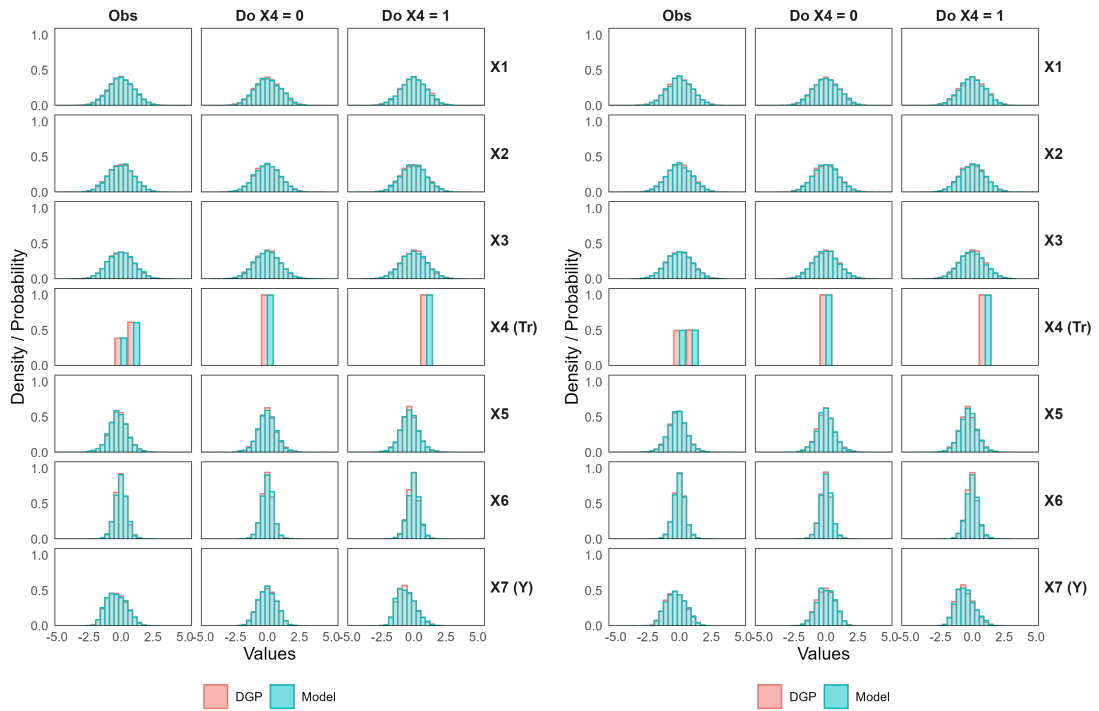


Figure 3.26: Marginal distributions of variables from the DGP and from samples generated by the fitted TRAM-DAG for Scenario 4.1 with direct and interaction effects. Distributions are shown as observed (Obs), under the control intervention $\text{do}(X_4 = 0)$, and under the treatment intervention $\text{do}(X_4 = 1)$. Left: Observational; Right: RCT setting.

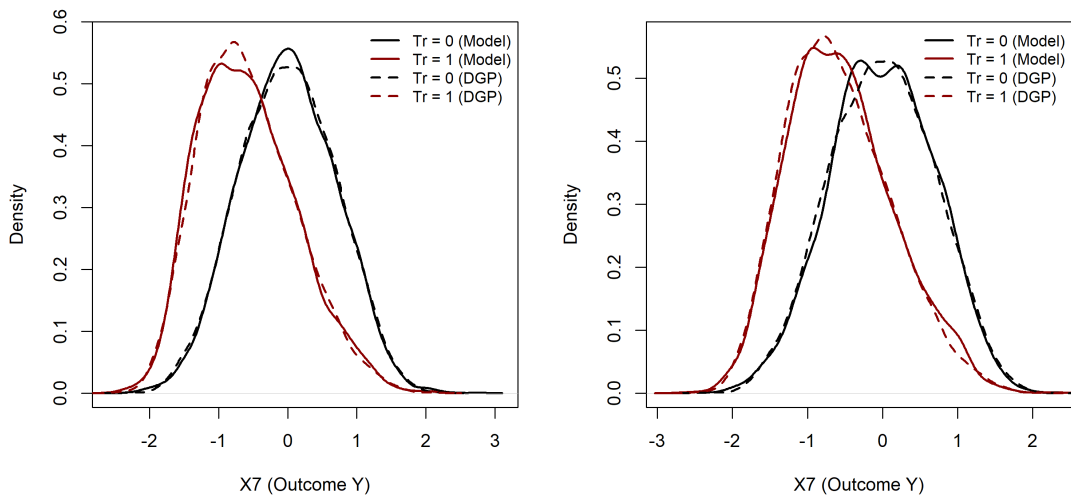


Figure 3.27: Distributions of the outcome variable ($X_7 = Y$) under control and treatment interventions for Scenario 4.1, which includes both direct and interaction effects. This plot provides a more detailed view of the X_7 panels shown under $\text{do}(X_4 = 0)$ and $\text{do}(X_4 = 1)$ in Figure 3.26. Left: Observational; Right: RCT setting.

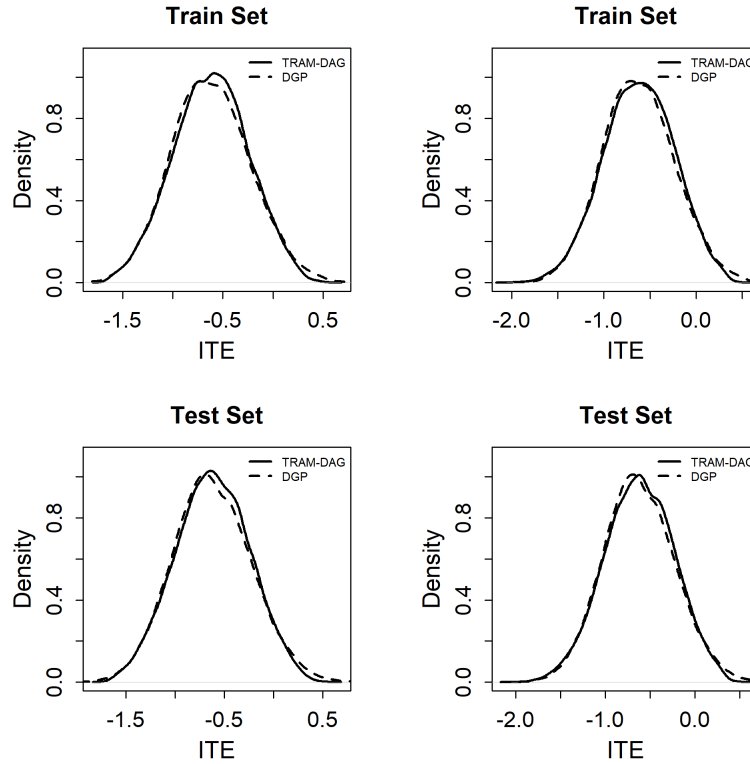


Figure 3.28: Densities of the estimated ITEs compared to the true ITEs in the training and test datasets for Scenario 4.1, which includes direct and interaction effects. Left: Observational; Right: RCT setting.

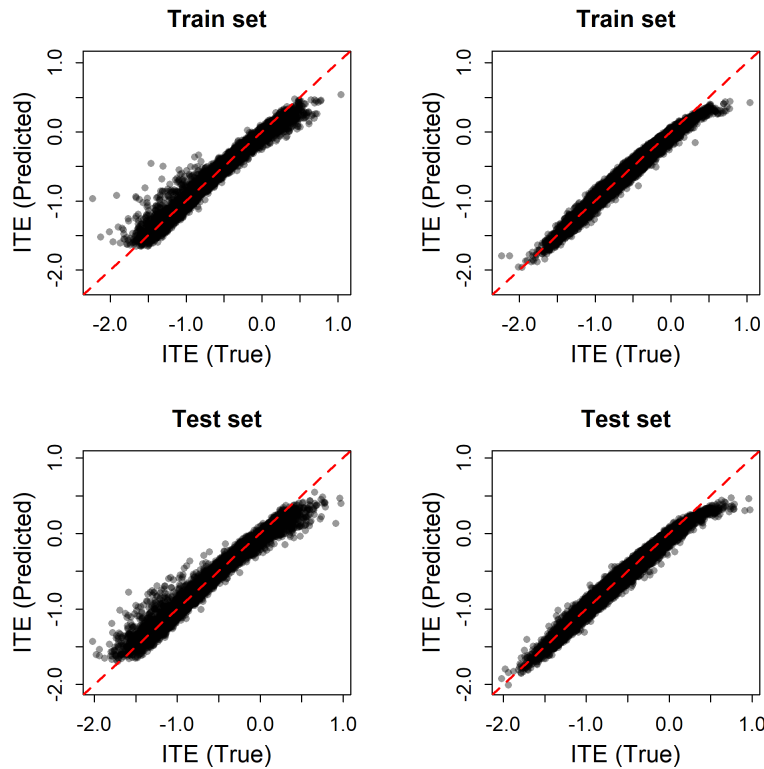


Figure 3.29: Scatterplots of estimated ITEs vs. true ITEs in the training and test datasets for Scenario 4.1, which includes direct and interaction effects. Left: Observational; Right: RCT setting.

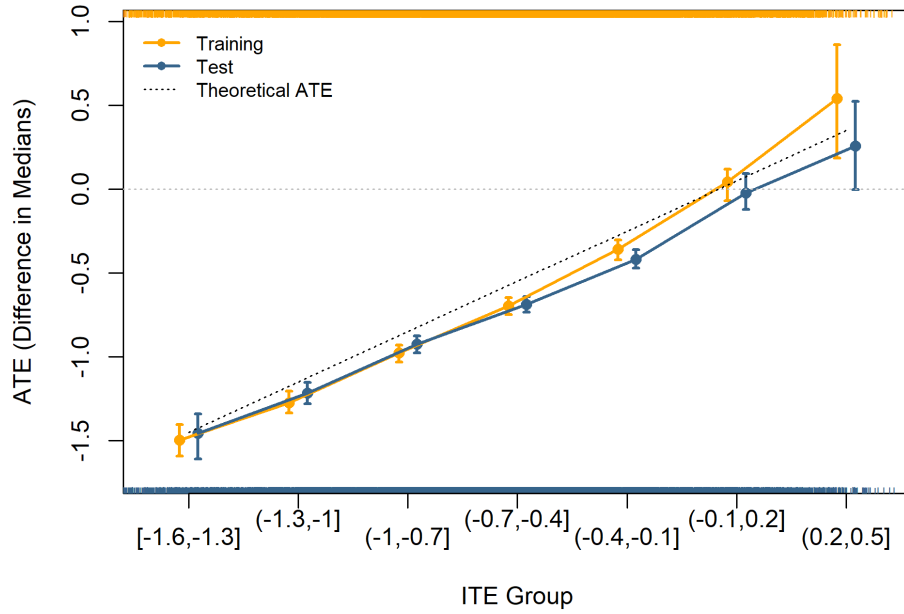


Figure 3.30: ITE-ATE plot for Scenario 4.1 in the observational setting, which includes direct and interaction effects. Individuals are grouped into bins based on their estimated ITEs, and within each bin, the ATE is calculated as the difference in medians of the observed outcomes under the two treatments. 95% bootstrap confidence intervals indicate the uncertainty.

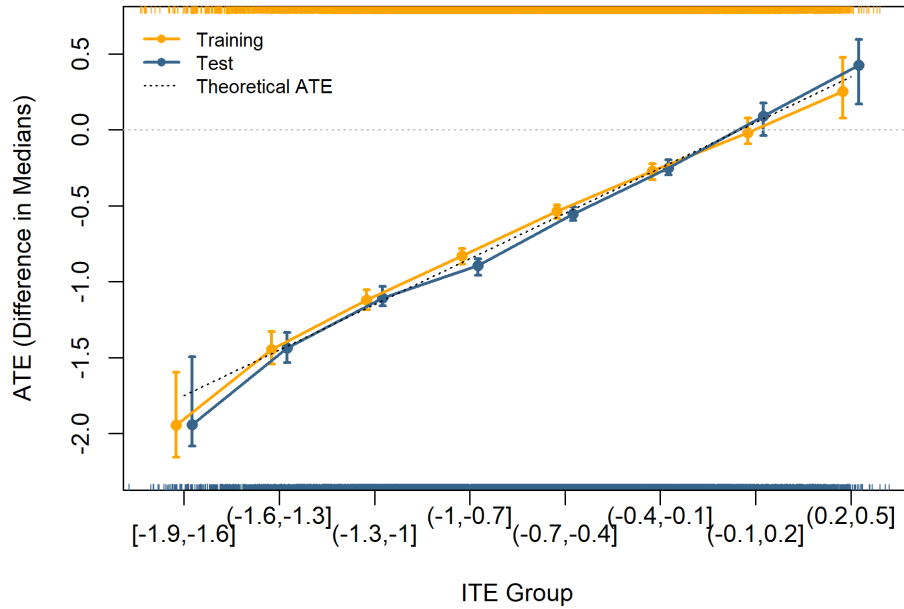


Figure 3.31: ITE-ATE plot for Scenario 4.1 in the RCT setting, which includes direct and interaction effects. Individuals are grouped into bins based on their estimated ITEs, and within each bin, the ATE is calculated as the difference in medians of the observed outcomes under the two treatments. 95% bootstrap confidence intervals indicate the uncertainty.

Discussion of Scenario 4.1: The TRAM-DAG successfully estimated the ITEs in Scenario 4.1 where a direct effect and interaction effects were present. There was no notable difference in estimation accuracy between the observational and RCT settings. Scatterplots of estimated ITEs vs. true ITEs showed good prediction accuracy (see Figure 3.29). The ATE based on the mean of the estimated ITEs closely matched the ATE derived from the true ITEs in both scenarios (see Table 3.1). These results highlight TRAM-DAG’s ability to compute counterfactuals for mediators and to estimate ITEs even in relatively complex DAG structures.

3.4.3.2 Scenario 4.2: With direct effect, but no interaction effects

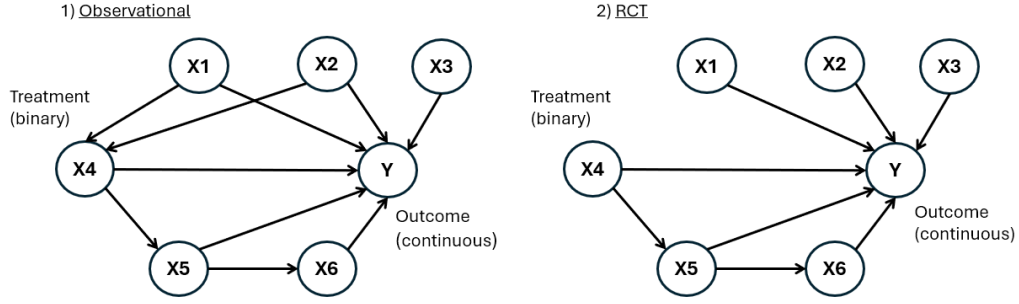


Figure 3.32: DAGs for Scenario 4.2, which includes a direct effect of the treatment on the outcome, but no interaction effects that would induce treatment effect heterogeneity. Left: Observational; Right: RCT setting.

Scenario 4.2 includes a direct effect of the treatment on the outcome, while the coefficients for the interaction terms are set to zero. This results in less heterogeneity in the ITE distribution compared to Scenario 4.1 (Section 3.4.3.1), as shown in Figure 3.33. The reason for the remaining heterogeneity despite the absence of explicit interactions is discussed at the end of this section. The observational and interventional densities generated by the fitted TRAM-DAG closely approximate the true DGP-defined distributions for all variables, as illustrated in Figures 3.34 and 3.35. However, there is a notable difference in variance between the estimated and true ITE distributions, visible in Figures 3.36 and 3.37. The ITE-ATE plots in Figures 3.38 and 3.39 are less informative than those in Scenario 1, which is as expected given the reduced heterogeneity. Table 3.2 presents the ATE measures for Scenario 2. In the test set of the RCT setting, the ATE based on the true ITEs was -0.633, while the ATE based on the estimated ITEs was -0.586.

In the RCT setting (training set), the difference in means of the outcomes between the two treatment groups was -0.569, with a 95% Wald confidence interval from -0.588 to -0.550. Note that the ATE in terms of difference in means cannot be directly compared to the ATE based on difference in medians.

Table 3.2: Scenario 4.2, including a direct treatment effect but no interaction effects: Comparison of ATE measures across train and test sets for the observational and RCT setting. $Y_{\text{observed}}^{(\text{Tr})}$ denotes the observed outcome under the treatment (Tr) actually received. Estimates based on these observed outcomes (means and medians) are provided only for the RCT setting, as the observational setting is confounded. The true ITEs (ITE_{true}) were calculated for each individual based on the data-generating process. In contrast, the estimated ITEs ($\text{ITE}_{\text{estimated}}$) were obtained from the TRAM-DAG trained on observed data. The estimated ATE from $\text{mean}(\text{ITE}_{\text{estimated}})$ can be directly compared to the true mean(ITE_{true}), whereas comparisons to empirical ATEs from observed outcome differences should be interpreted with caution. All ITEs were computed as quantile treatment effects (QTEs) based on the median of the potential outcome distributions, as defined in Equation 3.2.

Measure	Observational		RCT	
	Train	Test	Train	Test
ATE as $\text{mean}(Y_{\text{observed}}^{(1)}) - \text{mean}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.569	-0.572
ATE as $\text{median}(Y_{\text{observed}}^{(1)}) - \text{median}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.629	-0.639
ATE as $\text{mean}(\text{ITE}_{\text{true}})$	-0.633	-0.633	-0.633	-0.633
ATE as $\text{mean}(\text{ITE}_{\text{estimated}})$	-0.645	-0.644	-0.587	-0.586

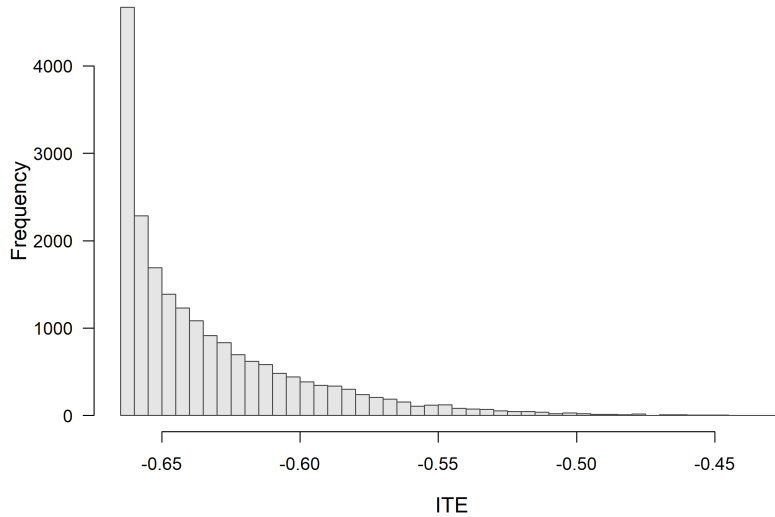


Figure 3.33: True ITE distribution resulting from the DGP for Scenario 4.2, which includes a direct treatment effect but no interaction effects. The true ITEs are identical in the observational and RCT settings, as they depend only on the potential outcomes under both treatment allocations.

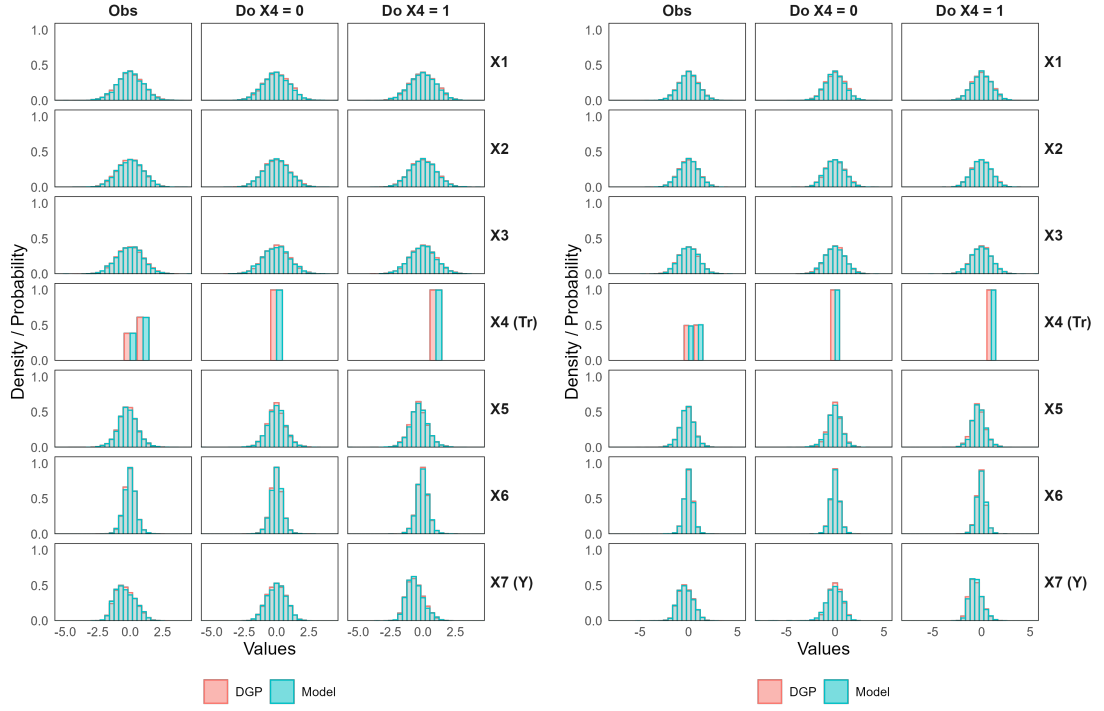


Figure 3.34: Marginal distributions of variables from the DGP and from samples generated by the fitted TRAM-DAG for Scenario 4.2, which includes a direct treatment effect but no interaction effects. The distributions are shown as observed (Obs), under control intervention $\text{do}(X_4 = 0)$, and under treatment intervention $\text{do}(X_4 = 1)$. Left: Observational; Right: RCT setting.

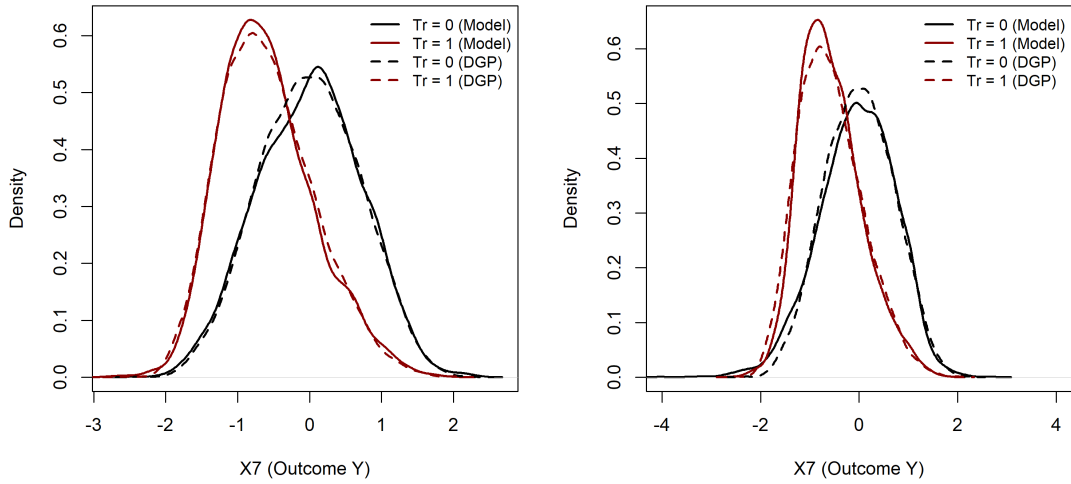


Figure 3.35: Distributions of the outcome variable ($X_7 = Y$) under treatment and control interventions for Scenario 4.2, which includes a direct treatment effect but no interaction effects. This plot provides a higher-resolution view of the X_7 panels under $\text{do}(X_4 = 0)$ and $\text{do}(X_4 = 1)$ from Figure 3.34. Left: Observational; Right: RCT setting.

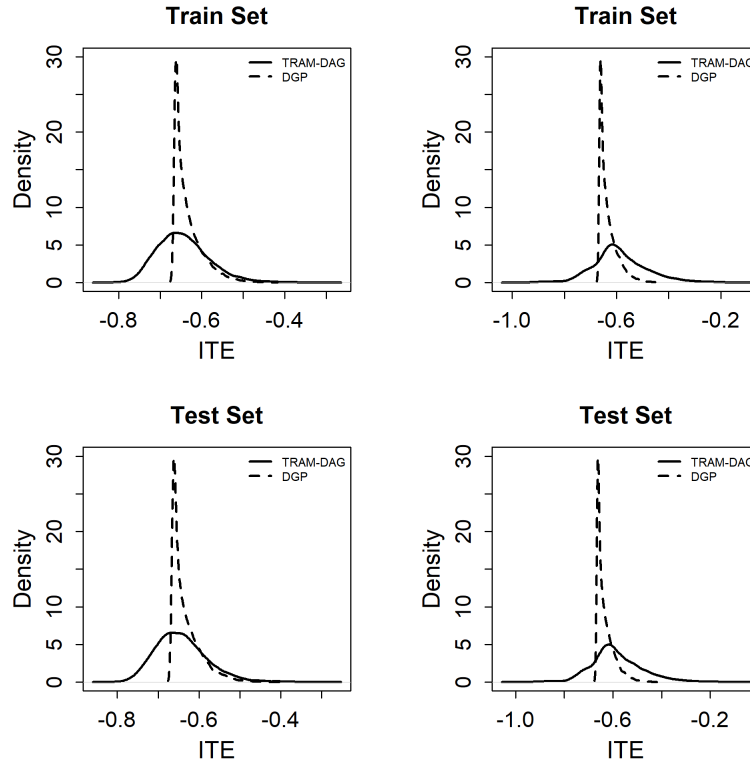


Figure 3.36: Densities of estimated ITEs compared to the true ITEs in the training and test datasets for Scenario 4.2, which includes a direct treatment effect but no interaction effects. Left: Observational; Right: RCT setting.

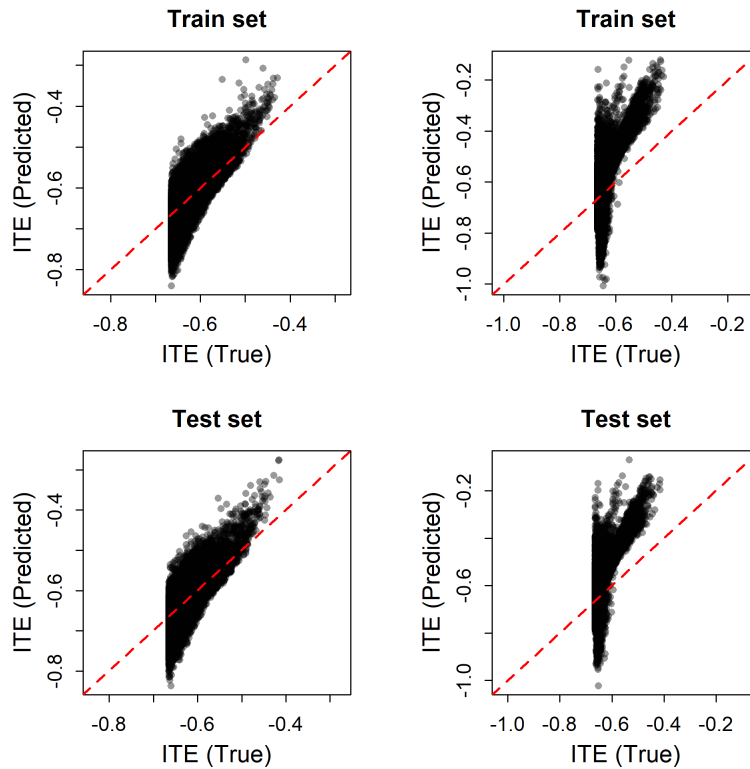


Figure 3.37: Scatterplots of estimated ITEs compared to the true ITEs in the training and test datasets for Scenario 4.2, which includes a direct treatment effect but no interaction effects. Left: Observational; Right: RCT setting.

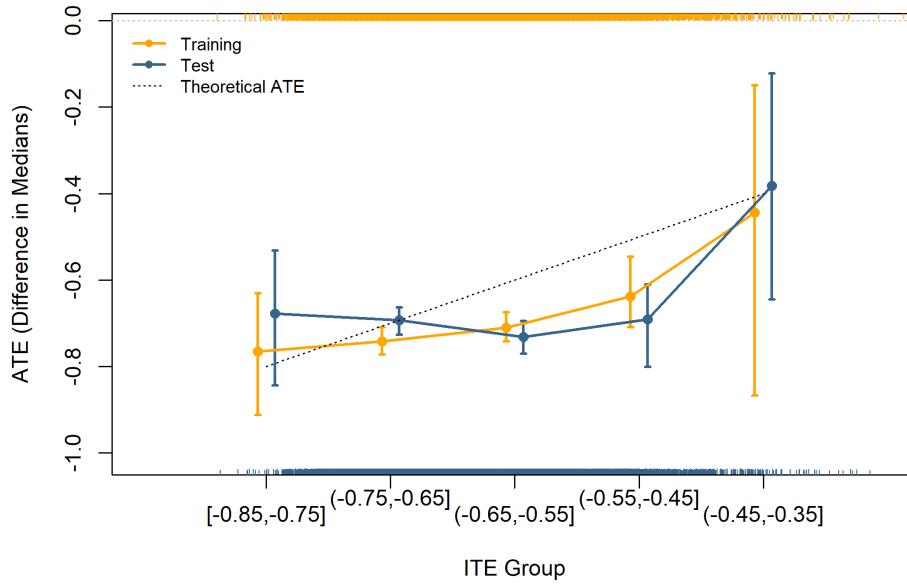


Figure 3.38: ITE-ATE plot for Scenario 4.2 in the observational setting, which includes a direct treatment effect but no interaction effects. Individuals are grouped into bins based on their estimated ITEs, and in each bin, the ATE is computed as the difference in medians of the observed outcomes under treatment and control. The 95% bootstrap confidence intervals indicate uncertainty.

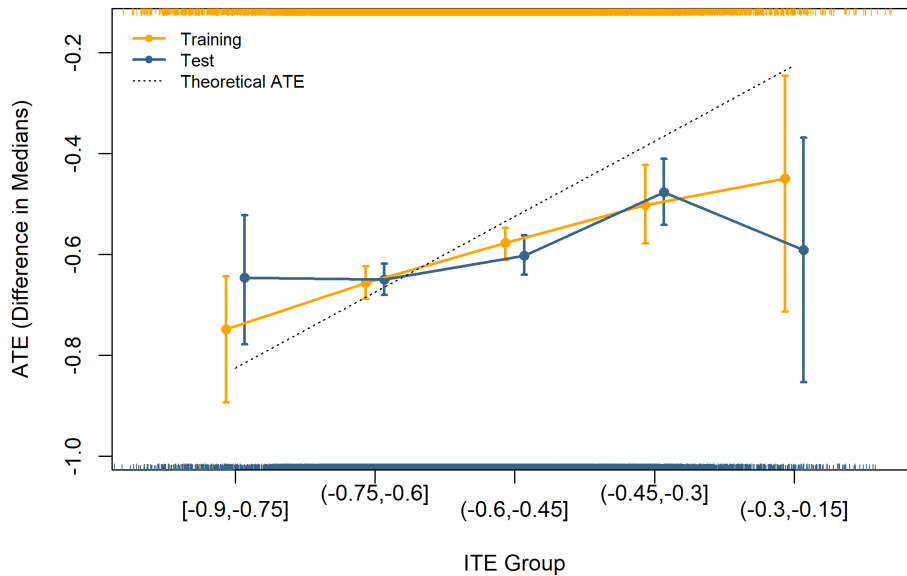


Figure 3.39: ITE-ATE plot for Scenario 4.2 in the RCT setting, which includes a direct treatment effect but no interaction effects. Individuals are grouped into bins based on their estimated ITEs, and in each bin, the ATE is computed as the difference in medians of the observed outcomes under treatment and control. The 95% bootstrap confidence intervals indicate uncertainty.

Discussion of Scenario 4.2: In Scenario 4.2, where no explicit interaction effects were present, ITE estimation was poor. Even though the TRAM-DAG could accurately reproduce observational and interventional distributions (see Figures 3.34 and 3.35), the resulting ITEs could not be accurately determined (see Figures 3.36 and 3.37). This aligns with our discovery in Experiment 3 (Section 3.3) that when true heterogeneity is weak, models tended to estimate too large heterogeneity, as shown, for example, in Figure 3.21 with the T-learner tuned random forest.

What might be surprising in Scenario 4.2 is the presence of heterogeneity (variation in true ITEs), despite the absence of explicitly specified treatment-covariate interactions in the data-generating process. As shown in Figure 3.33, one might expect constant ITEs across individuals – equal to the ATE – given the model’s additivity on the log-odds scale (Equation 3.1). However, as described by Hoogland *et al.* (2021), such heterogeneity arises because a constant treatment effect on the log-odds scale does not translate into a constant effect on a different scale, such as the probability scale. This phenomenon results from the nonlinearity of the inverse-link function (e.g., logit^{-1}), which transforms additive effects in the linear predictor into non-additive effects on the outcome scale. As they point out, the same shift induced by the treatment on the log-odds scale leads to different absolute risk reductions depending on the outcome risk under the control treatment. In other words, even with a homogeneous effect on the linear predictor, variation in covariates \mathbf{X} leads to different treatment effects on the probability scale.

This would not have occurred under a linear model where the transformation function h is the identity. In that case, the ITE would simplify as follows:

$$\text{ITE} = \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)] = (\beta_0 + \beta_t + \beta_x^\top \mathbf{X} + \epsilon) - (\beta_0 + \beta_x^\top \mathbf{X} + \epsilon) = \beta_t \quad (3.3)$$

In Equation 3.3, the ITE is constant and equal to the treatment coefficient β_t , independent of the covariates or the noise term, which both cancel out.

In contrast, under a nonlinear model, such as the logistic transformation model with a nonlinear intercept function used in this experiment, the ITE becomes:

$$\text{ITE} = \mathbb{E}[h^{-1}(Z + \beta_t + \beta_x^\top \mathbf{X})] - \mathbb{E}[h^{-1}(Z + \beta_x^\top \mathbf{X})] \quad (3.4)$$

Since h^{-1} is nonlinear, the difference depends on the covariate profile \mathbf{X} and on the noise term Z , even though the treatment effect β_t is additive in the linear predictor, i.e., on the log-odds scale. It may therefore be worth thinking about whether analyzing the ITE on a scale where the effect is constant offers any advantages.

3.4.3.3 Scenario 4.3: No direct effect, but with interaction effects

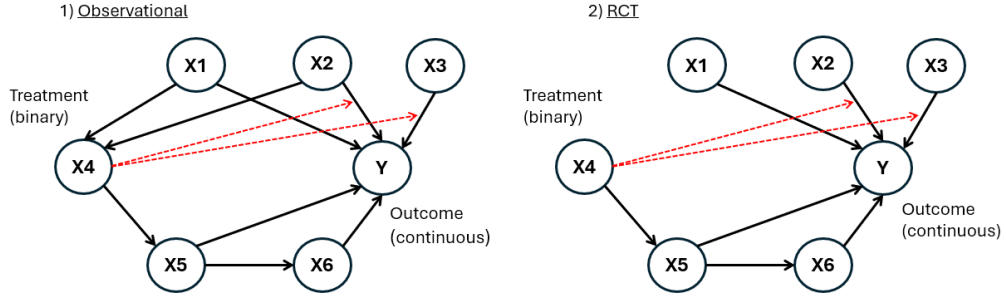


Figure 3.40: DAGs for Scenario 4.3, which includes no direct effect of the treatment on the outcome, but interaction effects with covariates X_2 and X_3 . Left: Observational; Right: RCT setting.

Scenario 4.3 includes no direct effect of the treatment on the outcome, but does include interaction effects between the treatment and covariates X_2 and X_3 . Compared to Scenario 4.1 (Section 3.4.3.1), removing the direct effect results in a more centered ITE distribution, as shown in Figure 3.41. In the test set of the RCT setting, the ATE measured as the difference in means was -0.048, with a 95% Wald confidence interval from -0.068 to -0.028. Note that the ATE in terms of difference in means cannot be directly compared to the ATE based on difference in medians.

Table 3.3: Scenario 4.3, without a direct treatment effect but including interaction effects: Comparison of ATE measures across train and test sets for the observational and RCT setting. $Y_{\text{observed}}^{(\text{Tr})}$ denotes the observed outcome under the treatment (Tr) actually received. Estimates based on these observed outcomes (means and medians) are provided only for the RCT setting, as the observational setting is confounded. The true ITEs (ITE_{true}) were calculated for each individual based on the data-generating process. In contrast, the estimated ITEs ($\text{ITE}_{\text{estimated}}$) were obtained from the TRAM-DAG trained on observed data. The estimated ATE from $\text{mean}(\text{ITE}_{\text{estimated}})$ can be directly compared to the true $\text{mean}(\text{ITE}_{\text{true}})$, whereas comparisons to empirical ATEs from observed outcome differences should be interpreted with caution. All ITEs were computed as quantile treatment effects (QTEs) based on the median of the potential outcome distributions, as defined in Equation 3.2.

Measure	Observational		RCT	
	Train	Test	Train	Test
ATE as $\text{mean}(Y_{\text{observed}}^{(1)}) - \text{mean}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.048	-0.048
ATE as $\text{median}(Y_{\text{observed}}^{(1)}) - \text{median}(Y_{\text{observed}}^{(0)})$	NA	NA	-0.048	-0.059
ATE as $\text{mean}(\text{ITE}_{\text{true}})$	-0.065	-0.068	-0.065	-0.068
ATE as $\text{mean}(\text{ITE}_{\text{estimated}})$	-0.059	-0.061	-0.051	-0.053

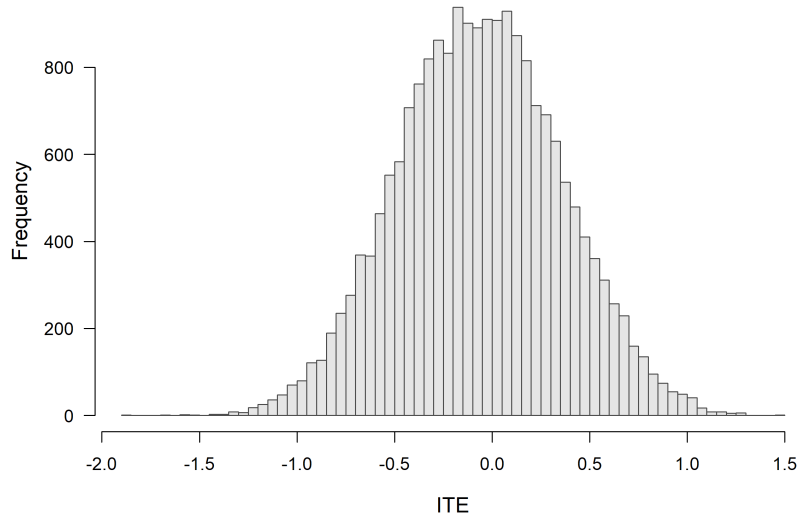


Figure 3.41: True ITE distribution resulting from the DGP for Scenario 4.3, which includes interaction effects but no direct treatment effect. The true ITEs are identical in the observational and RCT settings, since they are based on the potential outcomes under both treatment allocations.

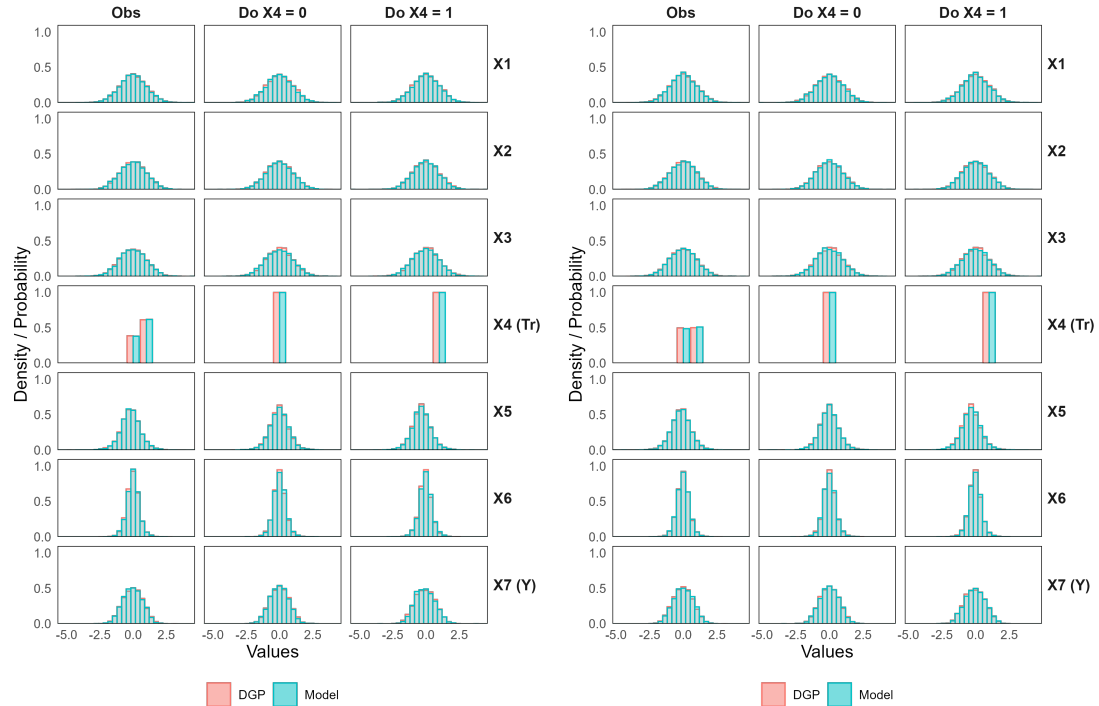


Figure 3.42: Marginal distributions of variables from the DGP and from samples generated by the fitted TRAM-DAG for Scenario 4.3, which includes interaction effects but no direct treatment effect. The distributions are shown as observed (Obs), under control intervention $\text{do}(X_4 = 0)$, and under treatment intervention $\text{do}(X_4 = 1)$. Left: Observational; Right: RCT setting.

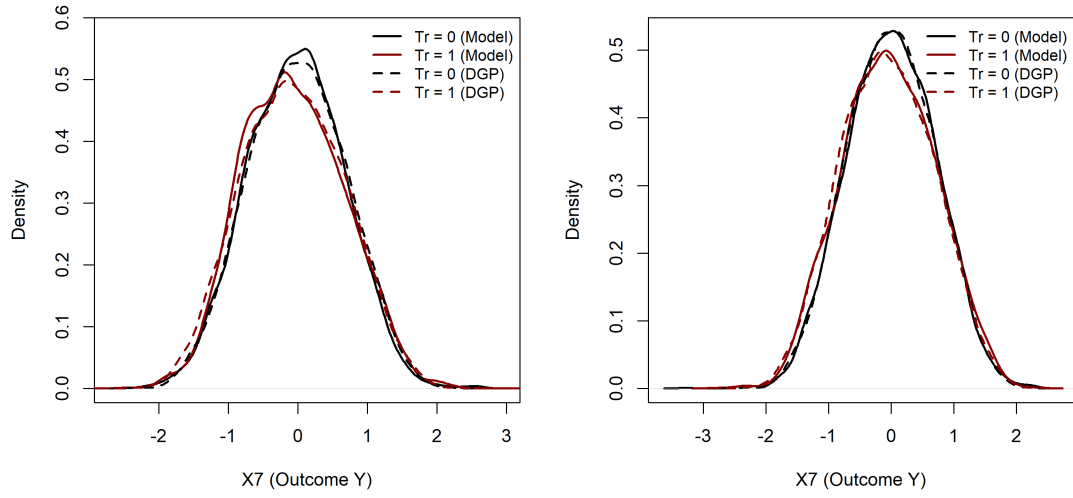


Figure 3.43: Distributions of the outcome variable ($X_7 = Y$) under treatment and control interventions for Scenario 4.3, which includes interaction effects but no direct treatment effect. This plot provides a higher resolution view of the X_7 panels under $\text{do}(X_4 = 0)$ and $\text{do}(X_4 = 1)$ from Figure 3.42. Left: Observational; Right: RCT setting.

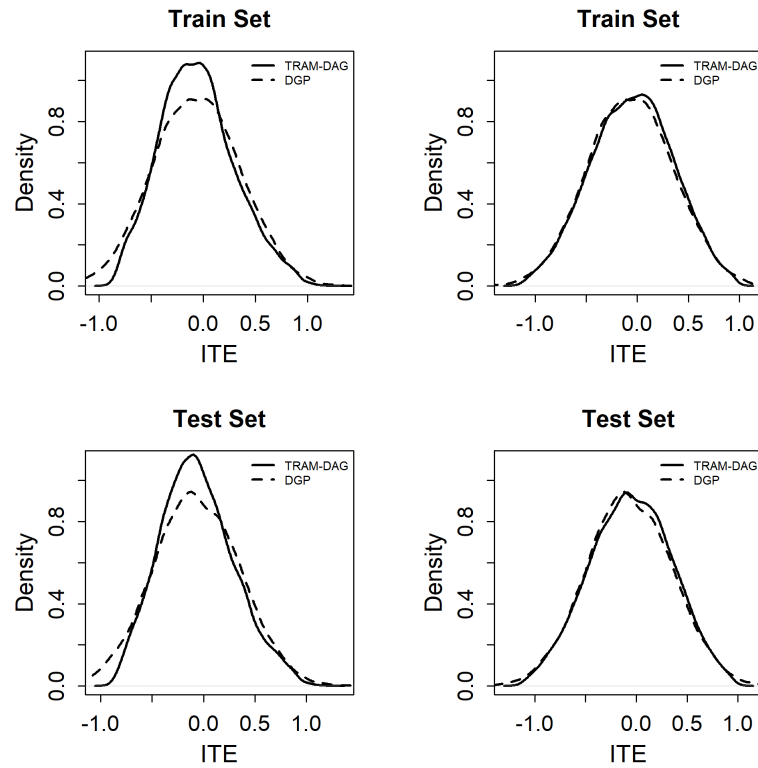


Figure 3.44: Densities of estimated ITEs compared to the true ITEs in the training and test datasets for Scenario 4.3, which includes interaction effects but no direct treatment effect. Left: Observational; Right: RCT setting.

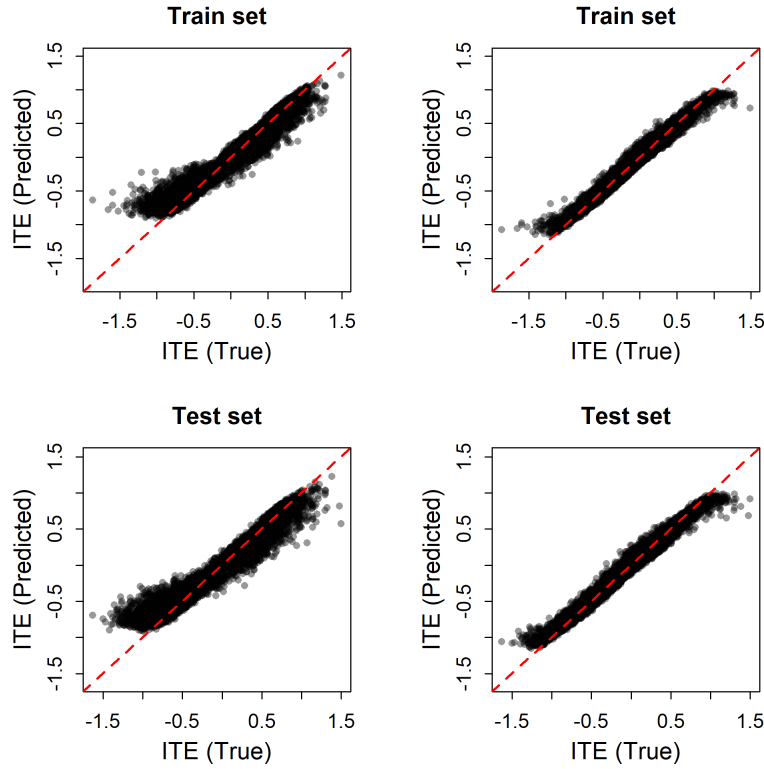


Figure 3.45: Scatterplots of estimated ITEs compared to the true ITEs in the training and test datasets for Scenario 4.3, which includes interaction effects but no direct treatment effect. Left: Observational; Right: RCT setting.

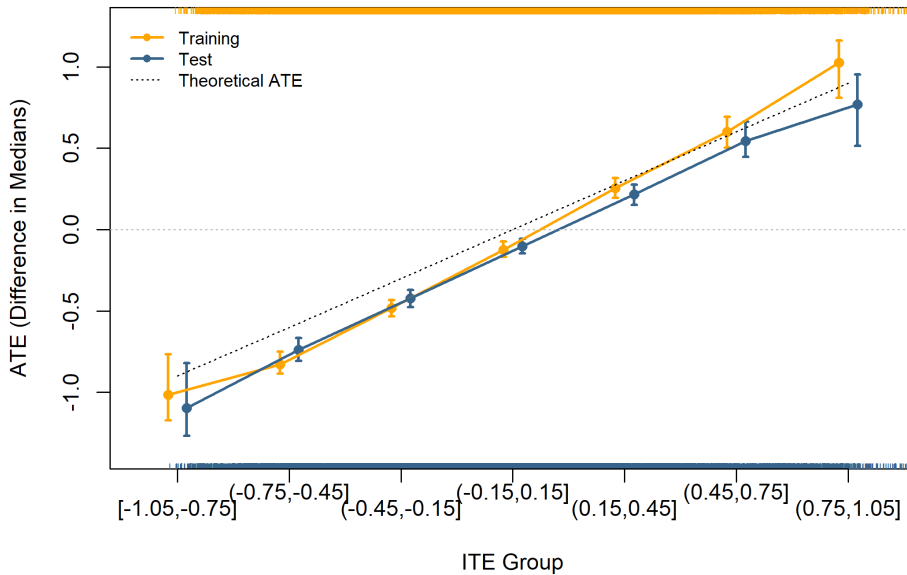


Figure 3.46: ITE-ATE plot for Scenario 4.3 in the observational setting, which includes interaction effects but no direct treatment effect. Individuals are grouped into bins based on the estimated ITE, and within each bin the ATE is computed as the difference in medians of the observed outcomes under treatment and control. 95% bootstrap confidence intervals reflect the uncertainty.

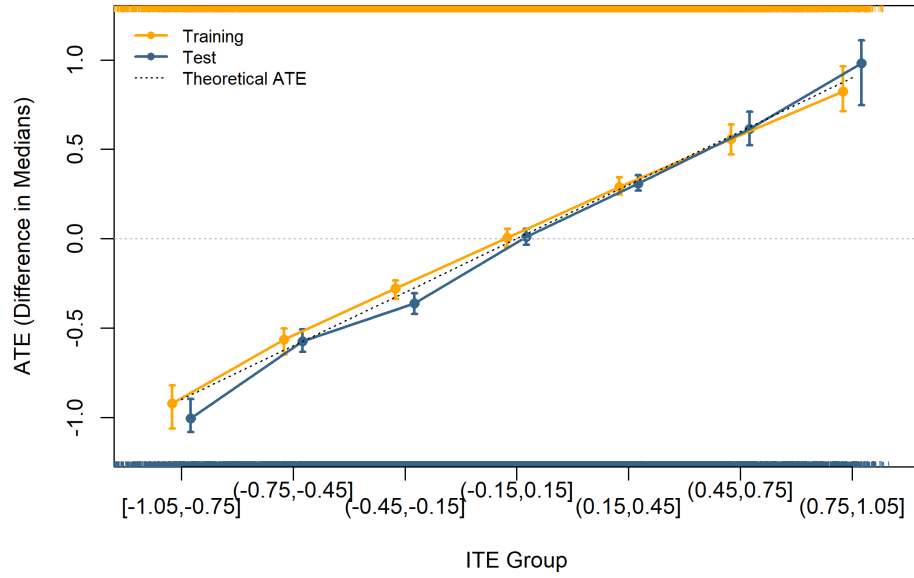


Figure 3.47: ITE-ATE plot for Scenario 4.3 in the RCT setting, which includes interaction effects but no direct treatment effect. Individuals are grouped into bins based on the estimated ITE, and within each bin the ATE is computed as the difference in medians of the observed outcomes under treatment and control. 95% bootstrap confidence intervals reflect the uncertainty.

Discussion of Scenario 4.3: In Scenario 4.3, which included interaction effects but no direct treatment effect, the TRAM-DAG could successfully estimate the ITE (similarly as in Scenario 4.1; Section 3.4.3.1). The scatterplots of estimated ITEs vs. true ITEs (Figure 3.45) showed good prediction accuracy.

3.4.4 Discussion of Experiment 4

TRAM-DAGs provided unbiased estimates of ITEs in Scenario 1 (Section 3.4.3.1) and Scenario 3 (Section 3.4.3.3), where heterogeneity of treatment effects was strong. This supports our claim that TRAM-DAGs can effectively be unbiased ITE estimation in a complex scenario, if the DAG is fully observed and heterogeneity is strong.

Chapter 4

Discussion

As the first focus of this thesis, we applied TRAM-DAGs as a flexible neural network-based approach to model causal relationships in a known directed acyclic graph (DAG). The second focus was on the estimation of individualized treatment effects (ITEs) in both randomized and confounded settings. We answered our research questions by conducting a dedicated experiment for each.

4.1 Findings

Experiment 1 (Chapter 3.1) served as a proof of concept, demonstrating that a TRAM-DAG can be fitted on simulated data from a known causal structure and used to sample from observational, interventional, and counterfactual distributions. Our results showed that the causal relationships could be recovered from the data, and that queries from all three levels of Pearl’s causal hierarchy could be accurately answered. We also gave advice on how to handle ordinal and categorical predictors (Appendix 5.4), how variable scaling affects the interpretation of coefficients (Appendix 5.5), and how interactions between variables can be modeled (Appendix 5.6).

Experiment 2 (Chapter 3.2) focused on the estimation of ITEs on the real dataset of the International Stroke Trial (IST). We tested whether we would reach the same conclusion as [Chen et al. \(2025\)](#) – that causal ML models fail to estimate ITEs that generalize to unseen data. We applied models based on logistic regression, random forest (with hyperparameter tuning), and TRAM-DAGs (with complex intercept to allow for interactions). The results showed that none of the three models produced ITE estimates that generalized to the independent test set. The observed treatment effect was not significant across any of the estimated ITE groups (i.e., when patients were grouped according to predicted ITEs). This motivated a deeper investigation into the possible reasons for poor model performance, which we addressed in Experiment 3.

Experiment 3 (Chapter 3.3) aimed to analyze the reasons for the poor ITE estimation performance observed in Experiment 2. We presented results for a logistic regression model (matching the DGP) and a tuned random forest, applied to ITE estimation in three simulated RCT scenarios. In the first scenario – the ideal case with a fully observed DAG and strong interaction effects – both models accurately recovered the true ITEs. However, as we demonstrated in Appendix 5.8, a default (untuned) random forest can perform poorly due to overfitting and miscalibration. In the second scenario, we found that omitting an important effect modifier led to biased ITE estimates; only a portion of the true heterogeneity was detected. In the third scenario, where true treatment effect heterogeneity was low, both models overestimated the heterogeneity. However, when validated on the test set, the ITEs did not generalize and no clear variation in treatment effects was observed. These results illustrate that when test set validation suggests no heterogeneity, it may be unclear whether the reason is a missing effect modifier or truly low treatment effect variation. Notably, in both cases, the identifiability assumptions are

not violated. Either of these scenarios could potentially explain the limited ITE performance on the IST dataset.

Experiment 4 (Chapter 3.4) aimed to demonstrate that TRAM-DAGs can be applied for ITE estimation in a confounded, relatively complex DAG, provided that the full DAG is observed. We concluded that TRAM-DAGs, with their ability to estimate causal relationships and model interactions between variables, yielded unbiased ITE estimates when interaction effects were present – in both confounded and randomized settings. In the case where no explicit treatment-covariate interaction was modeled in the DGP, we still observed some heterogeneity, which was attributable to the nonlinearity of the DGP and the TRAM-DAG model.

4.2 Limitations

Despite the promising results, this work has several limitations. While we aimed to make the simulation scenarios as realistic as possible while still retaining some interpretability, they may not fully reflect the complexity of real-world data. However, applying and evaluating models like TRAM-DAGs on real data for causal questions such as ITE estimation is inherently difficult, as the true effects are usually unknown. TRAM-DAGs rely on neural networks, which can be computationally demanding depending on model complexity and dataset size. And although TRAM-DAGs offer flexibility, we must still make modeling assumptions – for instance, regarding the scale of conditional effects – if interpretability is to be preserved.

4.3 Conclusion

Our findings showed that TRAM-DAGs were able to successfully recover causal relationships when the DAG was fully known and all relevant variables were observed. They also performed well for ITE estimation in controlled simulation settings. Through a series of experiments, we investigated the limitations of ITE estimation and concluded that poor performance may occur when important effect modifiers are unmeasured or when true treatment effect heterogeneity is low. These factors may help explain why ITE estimation failed to generalize in the real-world application on the International Stroke Trial dataset. We also observed that proper calibration of causal machine learning models is important to achieve accurate ITE estimates, though calibration alone may not be sufficient for valid individual-level predictions.

TRAM-DAGs, as generative causal models, allow for sampling from observational, interventional, and counterfactual distributions when fitted to a known DAG. Their ability to combine flexible neural network components with interpretable structure makes them well suited for real-world applications where both predictive accuracy and transparency are important.

Future work could apply TRAM-DAGs to additional real-world datasets, potentially including semi-structured data, to further explore the advantages of their modular design. It would also be valuable to investigate methods for improving ITE estimation in the presence of unobserved effect modifiers.

Overall, this thesis contributes to the growing field of causal inference in both observational and experimental settings, with a particular focus on individualized treatment effects and the capabilities of neural causal models.

Bibliography

- Allaire, J. and Tang, Y. (2025). *tensorflow: R Interface to 'TensorFlow'*. R package version 2.16.9. [14](#)
- Breiman, L. (2001). Random Forests. *Machine Learning*, **45**, 5–32. [25](#)
- Calster, B. V., van Smeden, M., Cock, B. D., and Steyerberg, E. W. (2020). Regression shrinkage methods for clinical prediction models do not guarantee improved performance: Simulation study. *Statistical Methods in Medical Research*, **29**, 3166–3178. [14](#)
- Chen, H., Aebersold, H., Puhane, M. A., and Serra-Burriel, M. (2025). Causal machine learning methods for estimating personalised treatment effects – insights on validity from two large trials. arXiv preprint 2501.04061. [iii](#), [v](#), [3](#), [4](#), [15](#), [20](#), [21](#), [23](#), [24](#), [55](#)
- Chernozhukov, V. and Hansen, C. (2005). An IV model of quantile treatment effects. *Econometrica*, **73**, 245–261. [33](#)
- Chollet, F., Allaire, J., et al. (2017). R Interface to 'Keras'. R package version 2.15.0. [14](#)
- Christensen, R., Bours, M. J., and Nielsen, S. M. (2021). Effect modifiers and statistical tests for interaction in randomized trials. *Journal of Clinical Epidemiology*, **134**, 174–177. [12](#)
- Elwert, F. and Winship, C. (2014). Endogenous selection bias: The problem of conditioning on a collider variable. *Annual Review of Sociology*, **40**, 31–53. [12](#)
- Freedman, B. (1987). Equipoise and the ethics of clinical research. *The New England journal of medicine*, **317**, 141–145. [1](#)
- Friedman, J. H., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, **33**, 1–22. [25](#)
- Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier neural networks. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, volume 15, 315–323. Proceedings of Machine Learning Research. [8](#)
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. (2017). On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70, 1321–1330. Proceedings of Machine Learning Research. [14](#)
- Gutierrez, P. and Gérardy, J.-Y. (2017). Causal inference and uplift modelling: A review of the literature. In *Proceedings of The 3rd International Conference on Predictive Applications and APIs*, volume 67, 1–13. Proceedings of Machine Learning Research. [1](#)
- Hariton, E. and Locascio, J. J. (2018). Randomised controlled trials - the gold standard for effectiveness research. *BJOG: An International Journal of Obstetrics & Gynaecology*, **125**, 1716 – 1716. [1](#)

- Hernan, M. and Robins, J. (2025). *Causal Inference: What If*. Chapman & Hall/CRC Press. CRC Press. [12](#)
- Herzog, L., Kook, L., Götschi, A., Petermann, K., Hänsel, M., Hamann, J., Dürr, O., Wegener, S., and Sick, B. (2023). Deep transformation models for functional outcome prediction after acute ischemic stroke. *Biometrical Journal*, **65**, 2100379. [8](#)
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association*, **81**, 945–960. [2](#), [12](#)
- Hoogland, J., Efthimiou, O., Nguyen, T. L., and Debray, T. P. A. (2024). Evaluating individualized treatment effect predictions: A model-based perspective on discrimination and calibration assessment. *Statistics in Medicine*, **43**, 4481–4498. [14](#)
- Hoogland, J., IntHout, J., Belias, M., Rovers, M. M., Riley, R. D., E. Harrell Jr, F., Moons, K. G. M., Debray, T. P. A., and Reitsma, J. B. (2021). A tutorial on individualized treatment effect prediction from randomized trials with a binary endpoint. *Statistics in Medicine*, **40**, 5961–5981. [12](#), [14](#), [24](#), [47](#)
- Hothorn, T., Kneib, T., and Bühlmann, P. (2014). Conditional transformation models. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, **76**, 3–27. [3](#), [5](#)
- Hothorn, T., Möst, L., and Bühlmann, P. (2018). Most likely transformations. *Scandinavian Journal of Statistics*, **45**, 110–134. [8](#), [18](#)
- International Stroke Trial Collaborative Group (1997). The International Stroke Trial (IST): a randomised trial of aspirin, subcutaneous heparin, both, or neither among 19,435 patients with acute ischaemic stroke. *The Lancet*, **349**, 1569–1581. [21](#)
- Ioffe, S. and Szegedy, C. (2015). Batch Normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, volume 37, 448–456. Proceedings of Machine Learning Research. [8](#)
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*. [8](#), [17](#)
- Kook, L. (2024). *comets: Covariance Measure Tests for Conditional Independence*. R package version 0.1-1. [21](#), [25](#)
- Kook, L., Herzog, L., Hothorn, T., Dürr, O., and Sick, B. (2022). Deep and interpretable regression models for ordinal outcomes. *Pattern Recognition*, **122**, 108263. [8](#)
- 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. [13](#)
- Nichols, A. (2007). Causal inference with observational data. *The Stata Journal*, **7**, 507 – 541. [1](#), [31](#)
- Papamakarios, G., Pavlakou, T., and Murray, I. (2017). Masked autoregressive flow for density estimation. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Proceedings of the 31st International Conference on Neural Information Processing Systems (NeurIPS)*, volume 30, 2335–2344. Curran Associates, Inc. [9](#)
- Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, **82**, 669–688. [10](#)

- Pearl, J. (2009a). Causal inference in statistics: An overview. *Statistics Surveys*, **3**, 96 – 146. [1](#)
- Pearl, J. (2009b). *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd edition. [2](#), [11](#)
- Poinsot, A., Leite, A., Chesneau, N., Sébag, M., and Schoenauer, M. (2024). Learning structural causal models through deep generative models: Methods, guarantees, and challenges. In *Proceedings of the 33rd International Joint Conference on Artificial Intelligence (IJCAI)*. [3](#)
- Prechelt, L. (2012). Early stopping — but when? In Orr, G. B. and Müller, K.-R., editors, *Neural Networks: Tricks of the Trade: Second Edition*, 53–67. Springer, Berlin, Heidelberg. [8](#)
- R Core Team (2024). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. [14](#)
- Reisach, A. G., Seiler, C., and Weichwald, S. (2021). Beware of the simulated DAG! Causal discovery benchmarks may be easy to game. In *Proceedings of the 35th Conference on Neural Information Processing Systems (NeurIPS)*. [65](#)
- Riley, R. D., Snell, K. I., Martin, G. P., Whittle, R., Archer, L., Sperrin, M., and Collins, G. S. (2021). Penalization and shrinkage methods produced unreliable clinical prediction models especially when sample size was small. *Journal of Clinical Epidemiology*, **132**, 88–96. [14](#)
- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, **70**, 41–55. [12](#)
- Rubin, D. B. (1980). Randomization analysis of experimental data: The Fisher randomization test comment. *Journal of the American Statistical Association*, **75**, 591–593. [12](#)
- Rubin, D. B. (2005). Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, **100**, 322–331. [12](#)
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, **323**, 533–536. [8](#)
- Sandercock, P. A., Niewada, M., Członkowska, A., and the International Stroke Trial Collaborative Group (2011). The International Stroke Trial database. *Trials*, **12**, 101. [21](#)
- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, **61**, 85–117. [8](#)
- Sick, B. and Dürr, O. (2025). Interpretable neural causal models with TRAM-DAGs. arXiv preprint 2503.16206, accepted at the CLear 2025 Conference. [iii](#), [3](#), [6](#), [9](#), [10](#), [11](#), [34](#)
- Sick, B., Hothorn, T., and Dürr, O. (2021). Deep transformation models: Tackling complex regression problems with neural network based transformation models. In *2020 25th International Conference on Pattern Recognition (ICPR)*, 2476–2481. [3](#), [5](#), [7](#)
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, **15**, 1929–1958. [8](#)
- Ushey, K., Allaire, J., and Tang, Y. (2025). *reticulate: Interface to 'Python'*. R package version 1.42.0. [14](#)
- Vegetabile, B. G. (2021). On the distinction between ”conditional average treatment effects” (CATE) and ”individual treatment effects” (ITE) under ignorability assumptions. arXiv preprint 2108.04939, presented at the Workshop on the Neglected Assumptions in Causal Inference (NACI), 38th International Conference on Machine Learning, 2021. [v](#), [31](#)

- Zhao, Z. and Harinen, T. (2020). Uplift modeling for multiple treatments with cost optimization. arXiv preprint 1908.05372. [1](#)
- Zheng, X., Aragam, B., Ravikumar, P., and Xing, E. P. (2018). DAGs with no tears: Continuous optimization for structure learning. arXiv preprint 1803.01422, accepted at NeurIPS 2018. [2](#)

Chapter 5

Appendix

5.1 Interpretation of linear coefficients

The transformation model framework allows for interpretation of the coefficients in the linear shift component. The choice of the inverse-link function F_Z determines the interpretation of the coefficients. For example, if the standard logistic distribution is chosen as the latent scale, i.e. $F_Z(z) = \text{expit}(z)$, the coefficients can be interpreted as log-odds ratios. For $F_Z(z) = 1 - \exp(-\exp(z))$, the interpretation of the coefficients changes to log-hazard ratios.

Consider the conditional transformation model:

$$F_{X_2|X_1}(x_2) = \text{expit}(h(x_2) + \beta_{12}x_1), \quad (5.1)$$

where $h(x_2)$ is a smooth, monotonic transformation (e.g., a Bernstein polynomial), and β_{12} is the coefficient representing the effect of X_1 on X_2 .

Applying the logit link (expit^{-1}) yields:

$$\log \left(\frac{F_{X_2|X_1}(x_2)}{1 - F_{X_2|X_1}(x_2)} \right) = h(x_2) + \beta_{12}x_1. \quad (5.2)$$

This shows that on the log-odds scale, X_1 has an additive linear effect on X_2 . The corresponding odds ratio when increasing x_1 by one unit is:

$$\text{OR}_{x_1 \rightarrow x_1+1} = \frac{\exp(h(x_2) + \beta_{12}(x_1 + 1))}{\exp(h(x_2) + \beta_{12}x_1)} = \exp(\beta_{12}). \quad (5.3)$$

Therefore, $\exp(\beta_{12})$ represents the multiplicative change in odds of $X_2 \leq x_2$ for a one-unit increase in X_1 , holding all else constant. This means that β_{12} can be interpreted as a log-odds ratio.

5.2 Bernstein polynomial for continuous outcomes

In deep TRAMs, the intercept for a continuous outcome y is modeled as a smooth and monotonically increasing function using a Bernstein polynomial of order M . Here, we focus on the more general case, where the intercept depends on the predictors \mathbf{x} . The same logic also applies to the simple intercept case without dependence on covariates.

The intercept of the transformation function can be written as:

$$h_I(y | \mathbf{x}) = \sum_{k=0}^M \vartheta_k(\mathbf{x}) \cdot B_{k,M}(s(y)), \quad (5.4)$$

where $s(y) \in [0, 1]$ is a scaled version of the outcome y , and $B_{k,M}(s(y))$ denotes the k -th Bernstein basis polynomial of degree M , defined as:

$$B_{k,M}(s(y)) = \binom{M}{k} s(y)^k (1 - s(y))^{M-k}.$$

The parameters $\vartheta_k(\mathbf{x})$ depend on the predictors and determine the shape of the intercept function. To ensure that $h_I(y | \mathbf{x})$ is monotonically increasing in y , the coefficients must satisfy:

$$\vartheta_0(\mathbf{x}) \leq \vartheta_1(\mathbf{x}) \leq \dots \leq \vartheta_M(\mathbf{x}).$$

To ensure monotonicity, the unbounded parameters $\tilde{\vartheta}_k(\mathbf{x}) \in \mathbb{R}$ are first predicted (by the neural network) and then transformed using a cumulative sum of softplus-transformed values. This guarantees that the resulting coefficients are non-decreasing, and therefore that the transformation function $h_I(y | \mathbf{x})$ is smooth and strictly monotonically increasing in y . Bernstein polynomials can approximate a wide range of smooth functions, provided the degree M is sufficiently large.

5.2.1 Scaling and extrapolation of the Bernstein polynomial

Because the Bernstein polynomial is only defined on the range $[0, 1]$, the outcome variable y must be scaled to the unit interval. For parameter estimation alone, this scaling would be sufficient. However, the transformation function $h(y | \mathbf{x})$ must also be evaluable at arbitrary values of y , particularly those outside the range of the training data. This is essential, for instance, when performing generative sampling, where predicted outcomes may lie beyond the originally observed range of y .

To address this, we extend the Bernstein polynomial with a linear extrapolation in the tails. Specifically, we construct the core transformation within the 5% and 95% quantiles of the training outcome y using the smooth Bernstein polynomial in Equation 5.4, and extrapolate beyond this range linearly using the boundary derivatives. This results in a piecewise-defined transformation that is smooth, differentiable, strictly monotonic, and defined for all real-valued outcomes.

Let $q_{0.05}$ and $q_{0.95}$ denote the 5th and 95th empirical quantiles of y , computed on the training set. The scaled outcome is defined as:

$$s(y) = \frac{y - q_{0.05}}{q_{0.95} - q_{0.05}}. \quad (5.5)$$

This transformation maps the central interval $[q_{0.05}, q_{0.95}]$ onto the unit interval $[0, 1]$, the domain of the Bernstein basis polynomials. Let $h_I(s(y) | \mathbf{x})$ be the core transformation function as defined in Equation (5.4). The extrapolated transformation $\tilde{h}_I(y | \mathbf{x})$ is then defined as:

$$\tilde{h}_I(y | \mathbf{x}) = \begin{cases} h_I(0 | \mathbf{x}) + h'_I(0 | \mathbf{x}) \cdot (s(y) - 0), & \text{if } s(y) < 0 \\ h_I(s(y) | \mathbf{x}), & \text{if } 0 \leq s(y) \leq 1 \\ h_I(1 | \mathbf{x}) + h'_I(1 | \mathbf{x}) \cdot (s(y) - 1), & \text{if } s(y) > 1 \end{cases} \quad (5.6)$$

The derivatives $h'_I(0 | \mathbf{x})$ and $h'_I(1 | \mathbf{x})$ are computed analytically from the Bernstein basis and the learned coefficients $\vartheta_k(\mathbf{x})$, ensuring smooth and differentiable transitions at the boundaries (see Section 5.2.2).

This construction ensures several desirable properties. First, the transformation function $\tilde{h}_I(y | \mathbf{x})$ is globally defined on \mathbb{R} , avoiding undefined regions or discontinuities. Second, monotonicity is guaranteed due to the softplus parameterization of the coefficients $\vartheta_k(\mathbf{x})$.

5.2.2 Analytical derivative of the Bernstein polynomial

To compute gradients and ensure differentiability at the extrapolation boundaries, we derive the analytical form of the derivative of the intercept of the transformation function.

Recall the general form of the intercept of the transformation function:

$$h_I(y \mid \mathbf{x}) = \sum_{k=0}^M \vartheta_k(\mathbf{x}) \cdot B_{k,M}(s(y)), \quad (5.7)$$

where $s(y) \in [0, 1]$ is the scaled outcome, $\vartheta_k(\mathbf{x})$ are the (monotonized) coefficients, and $B_{k,M}(s(y))$ are the Bernstein basis polynomials of degree M .

To compute the derivative with respect to y , we apply the chain rule:

$$\frac{d}{dy} h_I(y \mid \mathbf{x}) = \sum_{k=0}^M \vartheta_k(\mathbf{x}) \cdot \frac{d}{dy} B_{k,M}(s(y)) = \sum_{k=0}^M \vartheta_k(\mathbf{x}) \cdot \frac{dB_{k,M}(s)}{ds} \cdot \frac{ds}{dy}. \quad (5.8)$$

Since $s(y) = \frac{y - q_{0.05}}{q_{0.95} - q_{0.05}}$, its derivative is:

$$\frac{ds}{dy} = \frac{1}{q_{0.95} - q_{0.05}}. \quad (5.9)$$

The derivative of the Bernstein basis polynomial is:

$$\frac{d}{ds} B_{k,M}(s) = M [B_{k-1,M-1}(s) - B_{k,M-1}(s)]. \quad (5.10)$$

Therefore, the full derivative is:

$$\frac{d}{dy} h_I(y \mid \mathbf{x}) = \frac{M}{q_{0.95} - q_{0.05}} \sum_{k=0}^M \vartheta_k(\mathbf{x}) [B_{k-1,M-1}(s(y)) - B_{k,M-1}(s(y))]. \quad (5.11)$$

This expression is used to evaluate the slope of the transformation function at the borders and is also critical when computing the likelihood.

5.3 Negative log-likelihood

5.3.1 Continuous Outcome

For a continuous outcome Y , the conditional cumulative distribution function (CDF) is given by:

$$F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z(h(s(y) \mid \mathbf{x})), \quad (5.12)$$

where F_Z is the CDF of the standard logistic distribution:

$$F_Z(z) = \frac{1}{1 + e^{-z}}, \quad z \in \mathbb{R}, \quad (5.13)$$

and h is the conditional transformation function that maps the scaled outcome $s(y)$ to the latent (log-odds) scale.

The outcome y must be scaled to the unit interval $[0, 1]$ because the Bernstein polynomial is defined on this range:

$$s(y) = \frac{y - \min(y)}{\max(y) - \min(y)}. \quad (5.14)$$

To compute the conditional density, we apply the change-of-variables formula:

$$f_{Y|\mathbf{X}=\mathbf{x}}(y) = f_Z(h(s(y) | \mathbf{x})) \cdot h'(s(y) | \mathbf{x}) \cdot s'(y), \quad (5.15)$$

where f_Z is the PDF of the standard logistic distribution:

$$f_Z(z) = \frac{e^z}{(1 + e^z)^2}, \quad z \in \mathbb{R}. \quad (5.16)$$

The negative log-likelihood (NLL) contribution for a single observation is then given by:

$$\text{NLL} = -\log f_{Y|\mathbf{X}=\mathbf{x}}(y). \quad (5.17)$$

Plugging in the expressions yields:

$$\begin{aligned} \text{NLL} &= -\log f_{Y|\mathbf{X}=\mathbf{x}}(y) \\ &= -h(s(y) | \mathbf{x}) - 2 \log(1 + \exp(-h(s(y) | \mathbf{x}))) \\ &\quad + \log h'(s(y) | \mathbf{x}) - \log(\max(y) - \min(y)). \end{aligned} \quad (5.18)$$

5.3.2 Discrete Outcome

For a discrete outcome (binary, ordinal, categorical) with ordered categories y_k , $k = 1, \dots, K$, the transformation model defines the conditional CDF as:

$$F(Y \leq y_k | \mathbf{X} = \mathbf{x}) = F_Z(h(y_k | \mathbf{x})). \quad (5.19)$$

The likelihood contribution for an observation in class y_k is:

$$f_{Y|\mathbf{X}=\mathbf{x}}(y_k) = \begin{cases} F_Z(h(y_1 | \mathbf{x})), & k = 1, \\ F_Z(h(y_k | \mathbf{x})) - F_Z(h(y_{k-1} | \mathbf{x})), & k = 2, \dots, K-1, \\ 1 - F_Z(h(y_{K-1} | \mathbf{x})), & k = K. \end{cases} \quad (5.20)$$

The corresponding NLL contribution is then:

$$\text{NLL} = -\log f_{Y|\mathbf{X}=\mathbf{x}}(y_k). \quad (5.21)$$

5.4 Encoding of discrete variables

In TRAM-DAGs, a variable X_i can act either as a predictor variable for a child node or as an outcome variable that depends on parent nodes.

When X_i is the outcome variable, and it is discrete with K ordered categories (e.g., ordinal), its conditional distribution is modeled via a transformation function h that defines $K - 1$ cut-points. The modeling differences between continuous and discrete outcomes have already been discussed.

However, when a discrete variable X_i with K categories is used as a predictor variable, it should be dummy encoded. Dummy encoding creates $K - 1$ binary (0/1) indicator variables. Each binary variable corresponds to one of the non-reference categories, with the first category serving as the reference level that is not explicitly represented.

Example: Let X be an ordinal variable with three levels: 1, 2, 3. Dummy encoding results in two binary variables:

- X_1 : 1, if $X = 2$, 0 otherwise
- X_2 : 1, if $X = 3$, 0 otherwise

Now assume a continuous outcome Y that depends on X . The transformation model is:

$$F(Y | X) = F_Z(h_I(y) + x_1\beta_1 + x_2\beta_2)$$

This gives us the following cases:

- If $X = 1$ (reference level): $x_1 = 0, x_2 = 0$, so
 $F(Y | X = 1) = F_Z(h_I(y))$
- If $X = 2$: $x_1 = 1, x_2 = 0$, so
 $F(Y | X = 2) = F_Z(h_I(y) + \beta_1)$
- If $X = 3$: $x_1 = 0, x_2 = 1$, so
 $F(Y | X = 3) = F_Z(h_I(y) + \beta_2)$

The coefficients β_1 and β_2 represent the additive shift on the latent scale (e.g., log-odds) when moving from the reference category (1) to categories 2 and 3, respectively.

Dummy encoding ensures that discrete predictors can be incorporated into the deep TRAM framework and maintain interpretability.

5.5 Scaling of continuous variables

Neural networks work best when the input variables are standardized. While scaling (i.e., zero-mean and unit-variance normalization) removes the marginal variance pattern, [Reisach et al. \(2021\)](#) showed that many structure learning algorithms rely heavily on this pattern and may suffer a drop in performance when it is eliminated. However, since our approach assumes a known DAG structure, this is not a concern in our setting. Therefore, we scale continuous input variables in TRAM-DAGs.

A linear, monotonic, and invertible transformation of a predictor variable changes the interpretation of the coefficient. Scaling a predictor variable X as $X_{\text{std}} = (X - \mu_X)/\sigma_X$ implies that the coefficient $\tilde{\beta}$ is interpreted as the change in log-odds for a one standard deviation increase in the predictor variable – or equivalently, for a one unit increase in the standardized predictor. This differs from the interpretation of the original coefficient β , which represents the change in log-odds for a one-unit increase in the raw predictor variable.

In contrast, the standardization of the outcome variable does not affect the interpretation of the model, due to the scale invariance of the log-odds. Suppose we standardize the outcome Y as follows:

$$Y_{\text{std}} = \frac{Y - \mu_Y}{\sigma_Y}$$

This transformation is linear, monotonic, and invertible:

$$Y = Y_{\text{std}} \cdot \sigma_Y + \mu_Y$$

Therefore, for any threshold y , we have the equivalence:

$$P(Y < y | X) = P\left(Y_{\text{std}} < \frac{y - \mu_Y}{\sigma_Y} | X\right)$$

This means, the probability of being below a particular quantile remains the same after standardization. Consequently, the interpretation of coefficients in models with a continuous outcome remains unchanged. Specifically, the log-odds ratio

$$\log\left(\frac{P(Y < y | X + 1)}{1 - P(Y < y | X + 1)}\right) - \log\left(\frac{P(Y < y | X)}{1 - P(Y < y | X)}\right)$$

is equal to

$$\log \left(\frac{P \left(Y_{\text{std}} < \frac{y - \mu_Y}{\sigma_Y} \mid X + 1 \right)}{1 - P \left(Y_{\text{std}} < \frac{y - \mu_Y}{\sigma_Y} \mid X + 1 \right)} \right) - \log \left(\frac{P \left(Y_{\text{std}} < \frac{y - \mu_Y}{\sigma_Y} \mid X \right)}{1 - P \left(Y_{\text{std}} < \frac{y - \mu_Y}{\sigma_Y} \mid X \right)} \right)$$

as long as the same quantile (i.e., probability threshold) is used. Thus, the coefficient β reflects the same change in log-odds per one-unit increase in the (standardized) predictor, regardless of whether the outcome is standardized or not.

The general form of the transformation model is:

$$P(Y < y \mid X = x) = F_z(h(Y) + \beta \cdot X)$$

but now consider the case where this model is fitted using standardized outcome and predictors:

$$P(Y_{\text{std}} < y_{\text{std}} \mid X_{\text{std}} = x_{\text{std}}) = F_z(\tilde{h}(Y_{\text{std}}) + \tilde{\beta} \cdot X_{\text{std}})$$

where \tilde{h} and $\tilde{\beta}$ are the estimated transformation function and coefficients after standardizing the outcome and predictors.

Example: To evaluate the probability $P(Y < 20 \mid X = 3)$ in the standardized setting, we use:

$$P \left(\frac{Y - \mu_Y}{\sigma_Y} < \frac{20 - \mu_Y}{\sigma_Y} \mid X_{\text{std}} = \frac{3 - \mu_X}{\sigma_X} \right) = F_z \left(\tilde{h} \left(\frac{20 - \mu_Y}{\sigma_Y} \right) + \tilde{\beta} \cdot \frac{3 - \mu_X}{\sigma_X} \right)$$

In summary, standardizing predictors changes coefficient interpretation, whereas outcome standardization does not affect interpretability or model validity.

5.6 Modeling interactions with complex shift (CS)

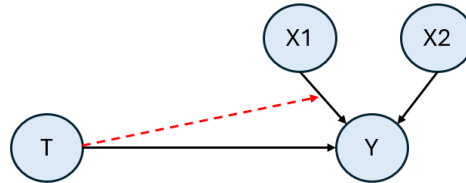


Figure 5.1: DAG for the complex shift example. The treatment T interacts with the covariate X_1 to influence the outcome Y . The covariate X_2 has a direct effect on the outcome.

The aim of this section is to show how TRAM-DAGs can effectively model interactions between multiple variables in a customizable way, using either a complex shift (CS) or a complex intercept (CI).

Here, we provide an example based on a simple DAG, as visualized in Figure 5.1. Assume that the binary treatment T is sampled from a Bernoulli distribution with probability 0.5. The covariates X_1 and X_2 represent baseline characteristics and are sampled from a standard normal distribution with a compound symmetric covariance matrix ($\rho = 0.1$). The outcome Y is binary and sampled from a Bernoulli distribution with a probability that depends on the treatment T and the covariates X_1 and X_2 . The treatment T interacts with the covariate X_1 to influence the outcome, while the covariate X_2 has a direct effect on the outcome. The formula for the outcome (on the log-odds scale) is given by:

$$\text{logit}(P(Y = 1 \mid T, X_1, X_2)) = \beta_0 + X_1\beta_1 + X_2\beta_2 + T\beta_3 + \textcolor{red}{TX_1}\beta_4$$

where $\beta_0 = 0.45$, $\beta_1 = -0.5$, $\beta_2 = 0.1$, $\beta_3 = -0.85$, and $\beta_4 = 0.7$. The interaction term TX_1 allows the effect of X_1 on the outcome to vary depending on the treatment group. This results in a heterogeneous treatment effect. We sample 20,000 training observations from this logistic data-generating process, and model it using a TRAM-DAG following the assumed structure in the DAG in Figure 5.1. The TRAM-DAG is specified with a complex shift (CS) of the treatment T and the covariate X_1 , which allows it to capture the interaction effect.

$$h(Y_k \mid T, X_1, X_2) = \vartheta_k + \textcolor{red}{CS}(T, X_1) + \text{LS}(X_2)$$

where $K = 1$ because the outcome is binary. The neural network for the CS is defined as shown in Figure 5.2, taking the treatment T and the covariate X_1 as inputs, which allows them to interact. After training the TRAM-DAG, we can visualize the resulting learned complex shift for each of the treatment groups, as illustrated in Figure 5.3. Due to the interaction, the effect of the covariate X_1 on the outcome Y differs between the treatment groups. The estimated effect aligns with the one specified in the data-generating process.

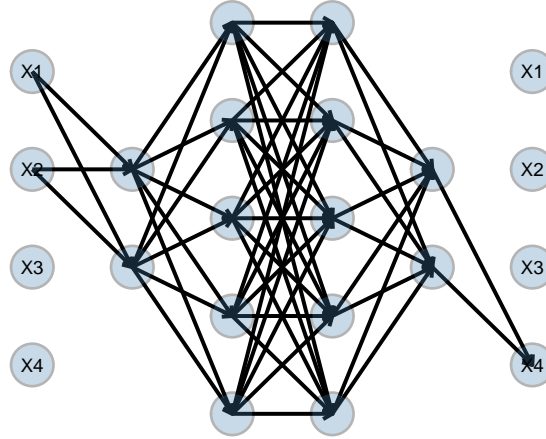


Figure 5.2: Neural network architecture for the complex shift (CS) effect of the treatment T and the covariate X_1 on the outcome Y (X_4). ReLU activation and batch normalization were used in the hidden layers. The deep neural network enables interactions between the input variables.

The fitted TRAM-DAG was then used to estimate the individualized treatment effect (ITE) for each sample i , defined as $\text{ITE}(\mathbf{x}_i) = P(Y_i = 1 \mid T = 1, \mathbf{X} = \mathbf{x}_i) - P(Y_i = 1 \mid T = 0, \mathbf{X} = \mathbf{x}_i)$. We estimated the ITEs on the test set using the trained TRAM-DAG. In the uncalibrated model, the predicted ITEs did not perfectly match the true ITEs – although the difference was only a small deviation. Therefore, we recalibrated the model by fitting a generalized additive model (GAM) on a separate calibration dataset using the predicted probabilities from the TRAM-DAG. The recalibrated ITEs aligned well with the true ITEs. The results are shown in Figure 5.4.

This example demonstrates how TRAM-DAGs can be used to model interactions, which are crucial for accurate ITE estimation when applying the model as an S-learner.

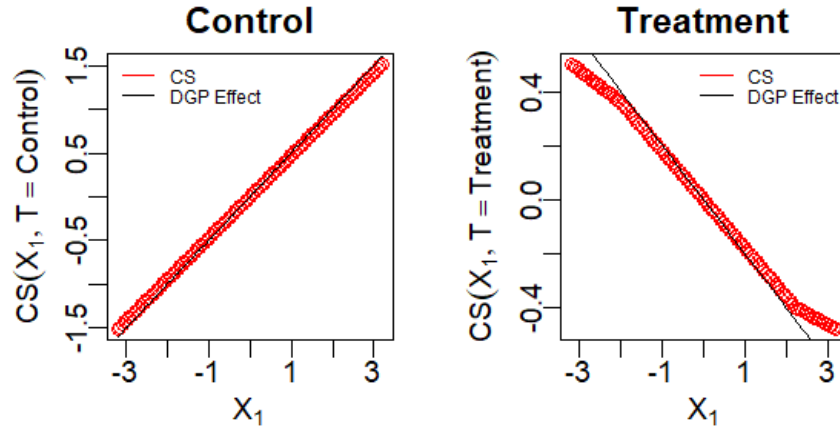


Figure 5.3: Learned complex shift (CS) effect of the treatment T and the covariate X_1 on the outcome (log-odds scale). Left: effect of X_1 on the outcome Y in the control group ($T = 0$); Right: effect of X_1 on the outcome Y in the treatment group ($T = 1$).

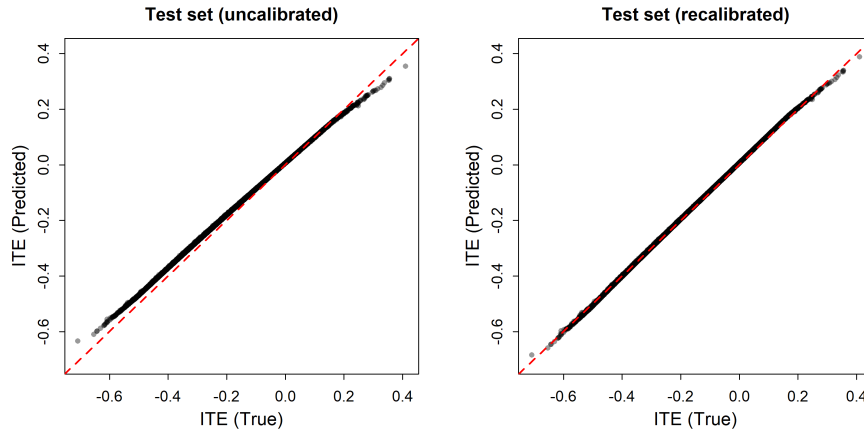


Figure 5.4: Estimated individualized treatment effects (ITE) using the fitted TRAM-DAG with a complex shift (CS). Left: estimated vs. true ITEs using the uncalibrated model. Right: estimated vs. true ITEs after recalibration. Recalibration was performed using an independent calibration dataset, where a generalized additive model (GAM) was fitted to recalibrate the probability estimates produced by the TRAM-DAG.

5.7 Experiment 2: Calibration plots

Figures 5.5-5.7 show the calibration plots of predicted risks versus observed outcome proportions for the models applied in Experiment 2 (International Stroke Trial, IST).

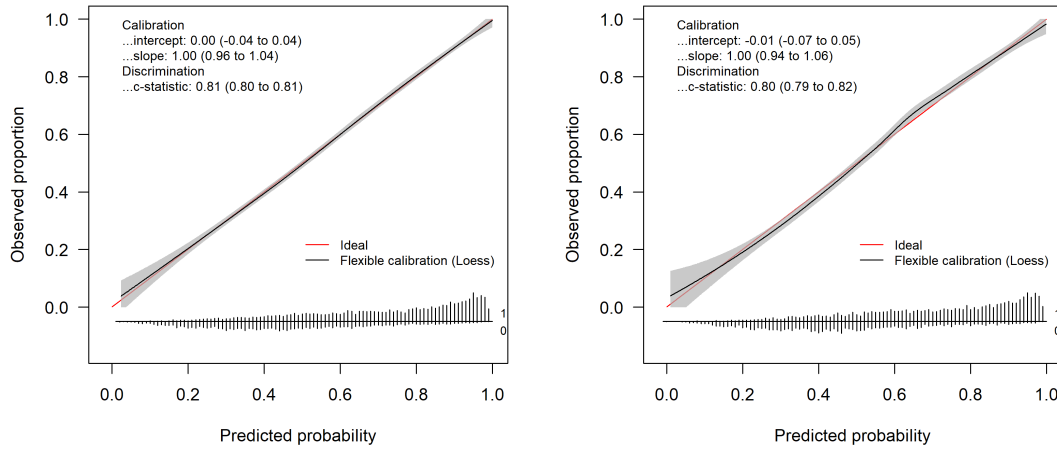


Figure 5.5: Calibration plot for the T-learner logistic regression applied to the International Stroke Trial (IST) in Experiment 2. The plot shows predicted risks versus observed event proportions. Left: training dataset; Right: test dataset.

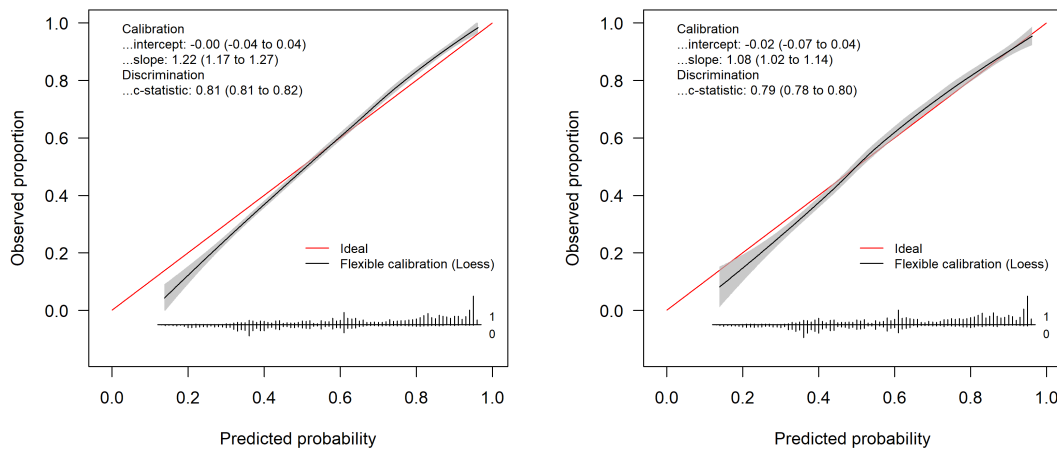


Figure 5.6: Calibration plot for the T-learner tuned random forest applied to the International Stroke Trial (IST) in Experiment 2. The plot shows predicted risks versus observed event proportions. Left: training dataset; Right: test dataset.

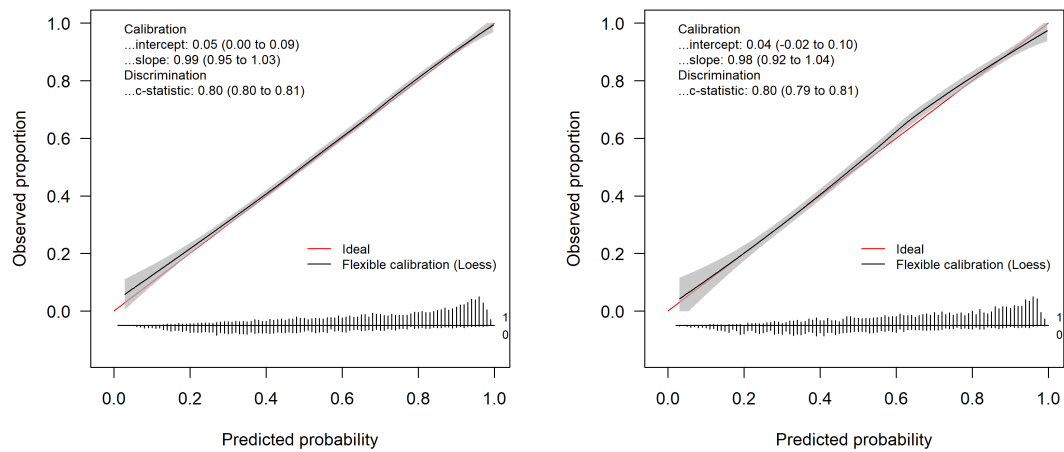


Figure 5.7: Calibration plot for the S-learner TRAM-DAG applied to the International Stroke Trial (IST) in Experiment 2. The plot shows predicted risks versus observed event proportions. Left: training dataset; Right: test dataset.

5.8 Experiment 3: Standard Random Forest for ITE Estimation

In Section 2.2, we emphasized the importance of model calibration when estimating individualized treatment effects (ITEs). Here, we present results from a default (untuned) random forest model in Scenario (1), where all variables are observed and both treatment and interaction effects are strong (see Figure 5.8). The corresponding model results are shown in Figure 5.9.

The scatterplot of true vs. predicted probabilities for $P(Y_i = 1 \mid \mathbf{X} = \mathbf{x}_i, T = t_i)$ in the training set shows that the model is poorly calibrated, with predicted probabilities deviating substantially from the true values. This miscalibration carries over to the ITE estimates. The ITE-ATE plot further illustrates that the estimated ITEs are not well aligned with the observed treatment effects, in both the training and test sets.

In contrast, the tuned random forest (Figure 3.15) achieves better calibration and more accurate ITE estimates. These results highlight the importance of proper model tuning for reliable ITE estimation, as poor calibration can lead to biased estimates.

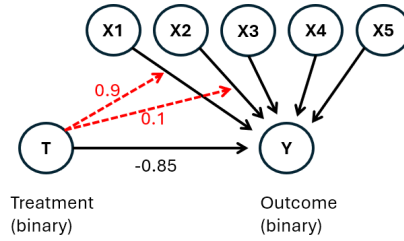


Figure 5.8: DAG for Scenario (1), where all variables are observed and both treatment and interaction effects are strong. This DAG was previously shown in Figure 3.12 and is re-plotted here for convenience. The numbers indicate the coefficients on the log-odds scale. Red arrows represent interaction effects between treatment (T) and covariates (X_1 and X_2) on the outcome (Y).

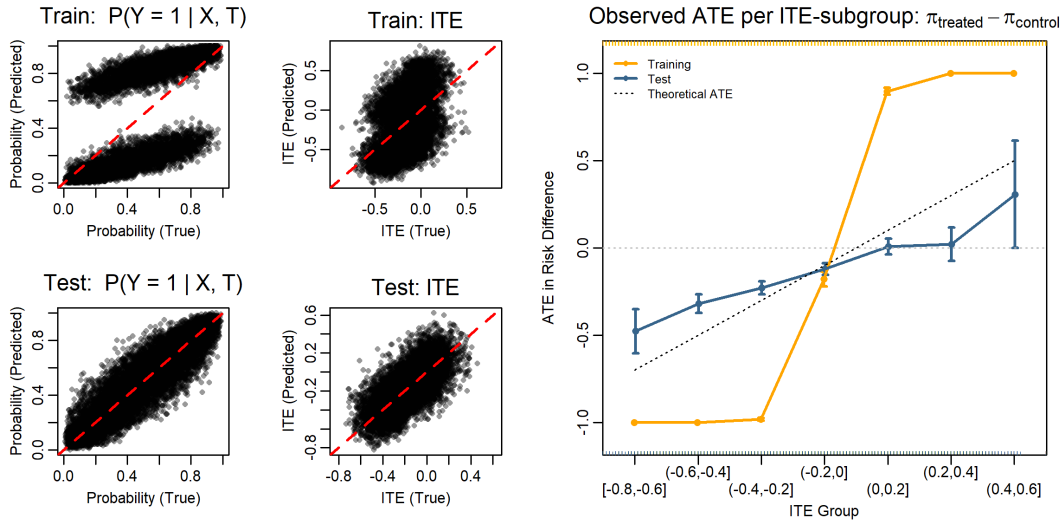


Figure 5.9: Results of the default random forest in Scenario (1), where the DAG is fully observed and both treatment and interaction effects are strong. Left: true vs. predicted probabilities for $P(Y = 1 \mid X, T)$; Middle: true vs. predicted ITEs; Right: observed ATE in terms of risk difference per estimated ITE subgroup.

5.9 Experiment 3: Calibration plots

Figure 5.10 shows the calibration plots in terms of the predicted risks against the the observed proportions of the event for the T-learner tuned random forest in Scenario (3.3) with weak direct and interaction effects. This is in contrast to the prediction plots presented in Figure 3.21 where we presented the true probabilities of the event $P(Y = 1 | X, T)$ against the predicted probabilities. It becomes apparent, that tuning the random forest model out-of-bag leads to a poor calibration on the training set (Figure 5.10; left), but due to improved generalization it leads to a better calibration on the test set (Figure 5.10; right).

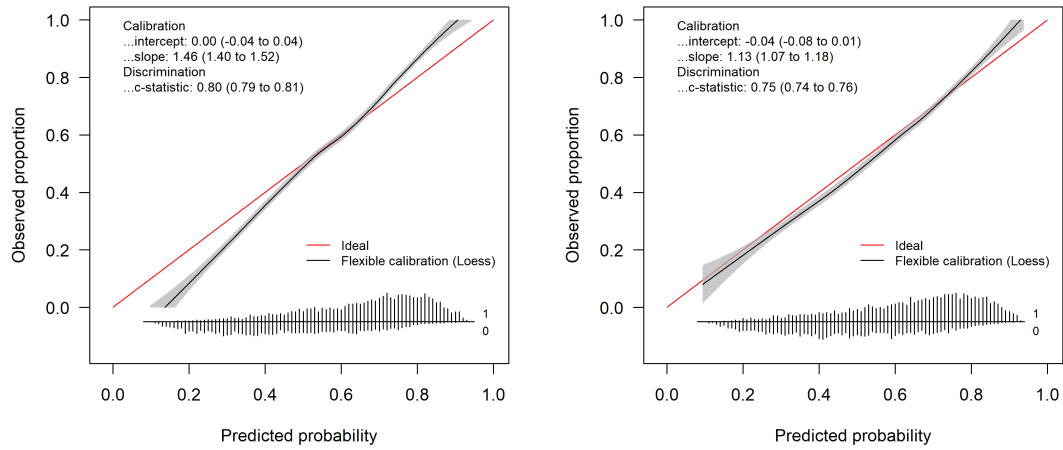


Figure 5.10: Calibration plot for the T-learner tuned random forest for Scenario (3) with weak direct and interaction treatment effects. It shows the predicted risks against the the observed proportions of the event. Left: training dataset; Right: test dataset.

5.10 Declaration of tools and services used

During the preparation of this thesis, I used ChatGPT, Google's Gemini and Github Copilot in order to support language refinement, such as checking grammar, spelling and clarity of expression as well as to assist in plotting and resolving some R coding errors. After using these tools/services, I reviewed and edited the content as needed and I take full responsibility for the content of the report.

