

CAUSALKANS: INTERPRETABLE TREATMENT EFFECT ESTIMATION WITH KOLMOGOROV-ARNOLD NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep neural networks achieve state-of-the-art performance in estimating heterogeneous treatment effects, but their opacity limits trust and adoption in sensitive domains such as medicine, economics and public policy. Building on well-established and high-performing causal neural architectures, we propose *causalKANs*, a framework that transforms neural estimators of *conditional average treatment effects* (CATEs) into Kolmogorov–Arnold Networks (KANs). By incorporating pruning and symbolic simplification, *causalKANs* yields interpretable closed-form formulas while preserving predictive accuracy. Experiments on benchmark datasets demonstrate that *causalKANs* perform on par with neural baselines in CATE error metrics, and that even simple KAN variants achieve competitive performance, offering a favorable accuracy–interpretability trade-off. By combining reliability with analytic accessibility, *causalKANs* provide auditable estimators supported by closed-form expressions and interpretable plots, enabling trustworthy individualized decision-making in high-stakes settings. We release the code for reproducibility.

1 INTRODUCTION

Estimating individual treatment effects from observational **tabular** data underpins high-stakes decisions in personalized medicine (Kent et al., 2018; Sanchez et al., 2022), public policy (Imai and Strauss, 2011), and economics (Manski, 2004), where interventions must be tailored beyond population averages (Wager and Athey, 2018). As personalized decision-making becomes the norm, accurately recovering conditional average treatment effects (CATEs) is indispensable for policy targeting and individualized care (Curth et al., 2024). However, accuracy alone is insufficient: regulatory frameworks (GDPR Art. 22; EU AI Act transparency for high-risk AI) and clinical practice increasingly discourage opaque models in consequential settings (European Parliament and Council, 2016; 2024; Goodman and Flaxman, 2017). Indeed, limited interpretability remains a barrier to the clinical adoption of machine learning (ML) systems (Amann et al., 2020; Tonekaboni et al., 2019). Yet contemporary state-of-the-art CATE estimators often rely on deep neural networks (Shalit et al., 2017; Shi et al., 2019), which we call *causalNN*, achieving strong performance but hindering auditing and trustable deployment.

We address this gap proposing *causalKANs*: a practical framework that transforms, or *KAN-ifies* (see Fig. 1), neural **tabular** CATE estimators into *closed-form*, auditable models by replacing their subnetworks with Kolmogorov–Arnold Networks (KANs) (Liu et al., 2024a;b). KANs parameterize one-dimensional edge functions (splines) and compose them via sums (and, in variants, products), enabling post-training simplification without departing from the trained predictor. Our interpretability notion is operational: after training, we apply edge-activity regularization, validation-guided pruning, and auto-symbolic substitution of learned splines with simple atoms to produce executable expressions for the potential outcomes and their difference (the CATE).

Concretely, our **contributions** are **threefold**. **First**, we introduce a **model-agnostic framework** for constructing KAN-based potential-outcome models from established causal neural architectures (e.g., metalearners (Künzel et al., 2019), TARNet (Shalit et al., 2017), DragonNet (Shi et al., 2019)) while

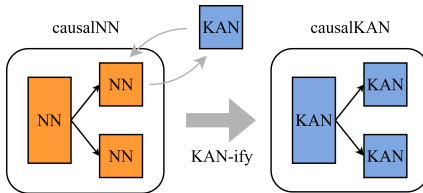


Figure 1: KAN-ification.

reusing their training objectives (see Fig. 1 for a sketch). **Second**, we propose a complete pipeline to achieve interpretability that includes pruning, symbolic substitution, and representation tools. **Third**, we present a comprehensive empirical study that i) compares performance of causalKANs on several well known benchmark datasets (IHDP (Hill, 2011), ACIC (Dorie et al., 2019), NSLM (Carvalho et al., 2019), NEWS (Johansson et al., 2016) and TCGA (Schwab et al., 2020)), and ii) assess the identification of causal equations with known synthetic datasets. The objective of benchmarking is to demonstrate that the use of KANs do not decrease the performance in causal metrics, namely, causalKANs are competitive with respect to causalNN, *not necessarily better*. For example, in Fig. 2 we can observe that we have not found statistical difference between T-KAN and DragonNet, which are the best models for IHDP A. In addition, the aim of ii) is to give practical reasons to think that if the true causal equations lie in the space that causalKANs can model, then the causal equation can be recovered up to an algebraic equivalence.

Our study positions causalKANs as domain-agnostic, accuracy-preserving, and inspection-ready CATE estimators: they retain the flexibility of deep architectures yet yield executable formulas that make effect modifiers and interactions explicit, aligning with emerging requirements for trustworthy, human-auditable decision support.

Empirically, causalKANs provide a favorable accuracy–interpretability frontier: simple heads (often additive or one hidden KAN layer) achieve *competitive* metrics across benchmarks, while additional depth rarely improves accuracy and consistently reduces formula compactness.

We present the preliminaries in §3, the causalKANs pipeline and instantiations in §4, and reports ablations, comparisons, and interpretability results in §5.

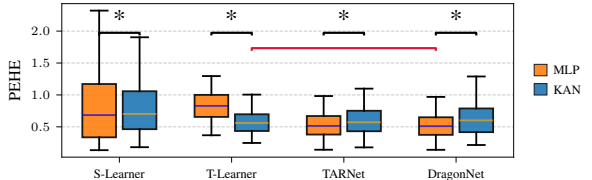


Figure 2: Overall, causalKANs achieves similar PEHE (lower is better) than causalNNs in IHDP A. * means statistical difference $p < 0.05$ in a Wilcoxon paired test. Red brace compares the best causalKAN with the best causalNN, without indicating statistical difference.

2 RELATED WORK

Kolmogorov-Arnold networks, as interpretable and parameter efficient alternatives of multi-layer perceptrons (MLP), have experienced a great growth in the last year, with the adaptation of existing technologies as convolutional networks (Bodner et al., 2024), residual networks (Yu et al., 2024), quantum networks (Wakaura et al., 2025), or autoencoders (Moradi et al., 2024). In the same manner, they have had an impressive adoption in domains where transparency is critical, as predicting risks in cardiovascular diseases (Al Bataineh et al., 2025) in healthcare (Almodóvar et al., 2025; Pendyala and Venkatachalam, 2025); predicting volatility in finances (Cho et al., 2025); predicting biomarkers (Alharbi et al., 2025), or predicting physics in power systems (Shuai and Li, 2025).

Interpretability is widely regarded as essential for the adoption of ML in sensitive domains (Doshi-Velez and Kim, 2017). Methods for improving interpretability based on KANs have been used in critical domains, such as Kolmogorov-Arnold Additive Models (KAAMs) (Almodóvar et al., 2025), which are shallow KANs that yields nonlinear additive models, which are well known as interpretable flexible estimators, as Neural Additive Models (NAMs) (Agarwal et al., 2021) or generalized additive models (Caruana et al., 2015). A closely related line is *symbolic regression* (SR), which searches over expression trees to recover closed-form models, classically via genetic programming (Schmidt and Lipson, 2009) and can produce highly concise formulas; however, SR is also computationally intensive and scales poorly (Zhang et al., 2022), while KANs, optimized by gradient descent, have been reported to be parameter-efficient and fast to optimize in practice (Liu et al., 2024b).

The need for interpretability is equally pressing in causal inference, where estimators are expected not only to predict counterfactuals but also to justify treatment decisions (Athey and Imbens, 2016). Early approaches such as causal trees (Athey and Imbens, 2016) provided transparent rule sets, and linear regressors were interpretable through their coefficients as causal parameters (Hahn et al., 2018). However, partially interpretable algorithms like causal forests (Foster et al., 2011) and BART (Hill, 2011) offered better predictive accuracy, while neural networks advanced performance even further (Johansson et al., 2016; Schwab et al., 2018; Shalit et al., 2017; Yao et al., 2018; Yoon et al., 2018) at

108 the cost of opacity. As a result, interpretability has often been relegated to a secondary concern in
 109 causal estimation (Rudin, 2018).

110 Bridging the trade-off between accuracy and transparency remains challenging. Some promising
 111 attempts include attention-based transformers highlighting confounder importance (Zhang et al.,
 112 2023), causal rule learners that yield decision rules but rely on black-box induction (Bargagli-Stoffi
 113 et al., 2020; Wu et al., 2023), fused lasso regression-based methods producing interpretable effect
 114 curves (Padilla et al., 2021), and, **importantly**, NAM-based estimators of average treatment effects
 115 (Chen et al., 2025). *Model distillation* approaches, such as training a surrogate on top of the work of
 116 Shalit et al. (2017) (Kim and Bastani, 2019), offer another path but depend on surrogate fidelity. Yet,
 117 addressing causal inference and interpretability together requires further work (Moraffah et al., 2020).
 118 In this context, KAN-inspired models such as our approach are promising because they combine the
 119 interpretability of additive structures with an expressive power undistinguishable of that of neural
 120 networks.

121 To the best of our knowledge, we are the second work that employs KANs for causal estimation, after
 122 KANITE (Mehendale et al., 2025), which replaces MLP backbones with KANs to estimate individual
 123 treatment effects under multiple treatments, using integral probability metrics/entropy-balancing
 124 variants and reporting gains in causal inference metrics over strong baselines. **Although this paper**
 125 **already performs a KAN-ification, it does not provide a systematic method for carrying it out, and**
 126 **its emphasis is predictive accuracy; it neither targets interpretability nor provides closed-form effect**
 127 **functions or visually interpretable plots. Instead, our work focuses on extracting human-readable**
 128 **causal effect formulas.**

129 Last, we also want to compare our proposal with the NAM-based approach of Chen et al. (2025). First,
 130 while NAMs offer decomposable architectures, they do not provide closed-form symbolic expressions
 131 nor a principled pipeline for transforming neural estimators into interpretable representations; second,
 132 the method proposed by Chen et al. (2025) is suitable only for ATE prediction, and third, our pipeline
 133 regularizes the functional behavior learned by neural networks, which can exhibit spiking or irregular
 134 patterns, while the simplification steps of our pipeline produce smoother, more stable mappings, which
 135 directly enhances interpretability and faithfulness in some real world problems (Liu et al., 2024b;
 136 Wang et al., 2025)

137 3 PRELIMINARIES

138 We consider an observational dataset $\mathcal{D} = \{\mathbf{x}_i, t_i, y_i\}_{i=1}^N$ of N i.i.d. samples from an unknown
 139 distribution $P(\mathbf{x}, \mathbf{t}, \mathbf{y})$, where $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^D$ are covariates, $\mathbf{t} \in \mathcal{T} = \{0, 1\}^1$ is the treatment, and
 140 $\mathbf{y} \in \mathbb{R}$ is the outcome. Our goal is to estimate the causal effect of the treatment on the outcome and
 141 provide a closed-form, interpretable formula for it.

142 Following the Neyman–Rubin potential outcomes framework (Rubin, 1974; Splawa-Neyman et al.,
 143 1990), each individual i has potential outcomes $\mathbf{y}_i(t)$ for $t \in \{0, 1\}$. Only the outcome corresponding
 144 to the received treatment is observed, which constitutes the *fundamental problem of causal inference*.
 145 The target estimand is the individual treatment effect (ITE): $\tau_i = \mathbf{y}_i(1) - \mathbf{y}_i(0)$.

146 Since τ_i is not observable, we estimate instead the conditional average treatment effect (CATE):

$$147 \tau(\mathbf{x}) := \mathbb{E}[\mathbf{y}(1) \mid \mathbf{x} = \mathbf{x}] - \mathbb{E}[\mathbf{y}(0) \mid \mathbf{x} = \mathbf{x}], \quad \text{denoting} \quad \mu_t(\mathbf{x}) := \mathbb{E}[\mathbf{y}(t) \mid \mathbf{x} = \mathbf{x}]. \quad (1)$$

148 We assume the standard conditions of causal inference: **i) positivity**, $0 < P(\mathbf{t} = t \mid \mathbf{x}) < 1$; **ii)**
 149 **conditional ignorability**, $\mathbf{y}(t) \perp\!\!\!\perp \mathbf{t} \mid \mathbf{x}$, which requires a valid adjustment set blocking all backdoor
 150 paths²; **iii) consistency**, $\mathbf{y}_i(t_i) = y_i$; and **iv) no interference**, $\mathbf{y}_i \perp\!\!\!\perp \mathbf{y}_j$ for $i \neq j$. Under these
 151 assumptions, it follows that $\mathbb{E}[\mathbf{y} \mid \mathbf{x}, \mathbf{t}]$ equals $\mu_t(\mathbf{x})$ (Hernán and Robins, 2025). The challenge is
 152 that this conditional expectation becomes increasingly difficult to estimate with high-dimensional,
 153 multimodal covariates and complex covariate–treatment–outcome relations, motivating the use of
 154 flexible function approximators such as neural networks.

155 ¹We focus on binary treatments for clarity, though the derivations extend to multi-valued and in some cases
 156 continuous treatments; see App. A.1.

157 ²An adjustment set must block all backdoor paths between treatment and outcome, avoid selection bias, and
 158 not block front-door paths.

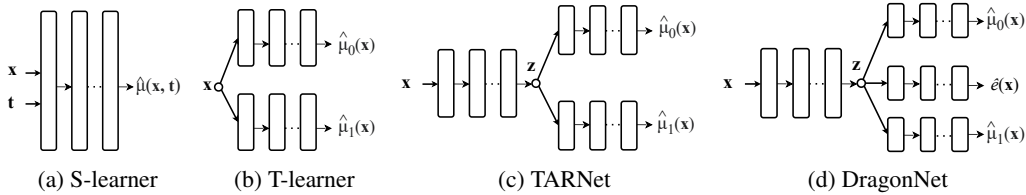


Figure 3: Architectures used for *potential outcome regression*. Boxes denote layers or backbones (neural networks or KANs); dotted arrows indicate optional hidden layers.

3.1 CAUSAL NEURAL NETWORKS

Neural networks are flexible function approximators (Hornik et al., 1989) and often provide state-of-the-art potential–outcome and CATE estimates (Alaa and Schaar, 2018; Curth and Van der Schaar, 2021b; Tesei et al., 2023). We summarize the canonical architectures used in our experiments (see Fig. 3); all perform potential–outcome regression, $\mu_t(x)$, and obtain the CATE by difference, $\hat{\tau}(x) = \hat{\mu}_1(x) - \hat{\mu}_0(x)$. Our replacement of neural backbones by KANs is orthogonal to these designs and can also be applied to methods that target CATE directly (e.g., X-, R-learners, MRIV-Net; Frauen and Feuerriegel, 2022; Künzel et al., 2019; Nie and Wager, 2021).

Meta-learners (Künzel et al., 2019) are model-agnostic. The *S-learner* fits a single regressor $\mu_t(x, t)$ of the factual outcome and obtains potential outcomes by toggling t . Its simplicity and ability to handle continuous treatments are appealing, but it may underuse t when the treatment signal is weak relative to x . The *T-learner* trains two separate regressors, $\hat{\mu}_0(x)$ and $\hat{\mu}_1(x)$, one per arm. This allows treatment-specific fits but splits the data, which can increase variance and yield sharp CATE estimates in small samples.

TARNet (Shalit et al., 2017) introduces a shared representation $z(x)$ feeding two heads, combining the strengths above: shared structure across arms (mitigating data inefficiency) with treatment-specific outcome mappings (reducing treatment ignorance). Its CFRNet variant changes only the loss to encourage balanced representations, while the architecture remains identical. TARNet has shown robust performance and is widely used as a baseline (Curth and Van der Schaar, 2021a;b; Schwab et al., 2018; Shalit et al., 2017; Yao et al., 2018).

DragonNet (Shi et al., 2019) extends TARNet with a third head for the propensity score, $\hat{e}(x)$, encouraging $z(x)$ to capture confounding structure through multitask learning. Originally paired with targeted regularization for doubly robust average treatment effect (ATE) estimation (Van der Laan et al., 2011), the base network already yields accurate potential outcomes and thus CATEs (Curth and Van der Schaar, 2021b; Ling et al., 2023; Lolak et al., 2025). Intuitively, predicting both outcomes and treatment forces the representation to retain covariates that co-determine assignment and response, which benefits counterfactual prediction.

Fig. 3 illustrates the corresponding schematics; we keep the architectures unchanged when replacing neural components by KANs.

3.2 KOLMOGOROV–ARNOLD NETWORKS

We employ KANs as our backbones. Kolmogorov–Arnold Networks (KANs) (Liu et al., 2024b) are deep models that replace fixed node-wise activations with *learnable univariate functions on edges*, using *addition*—and, in extensions, multiplication (Liu et al., 2024a)—as the only explicit multivariate operations. Their design is motivated by the Kolmogorov–Arnold representation theorem (KART), which guarantees that any continuous $f : \mathbb{R}^D \rightarrow \mathbb{R}$ can be expressed as

$$f(x) = \sum_{q=1}^{2D+1} \Phi_q \left(\sum_{k=1}^D \phi_{q,k}(x_k) \right), \quad (2)$$

where all nonlinearities are one-dimensional (Arnold, 1957; Braun and Griebel, 2009; Kolmogorov, 1956; 1957). This motivates stacked architectures where multivariate structure emerges from compositions of simple univariate transformations and aggregations. Formally, a depth- L KAN with widths

$\{n_l\}$ and coordinates $z_{l,r}$ replaces the linear map and fixed activation at layer l by

$$z_{l+1,s} = \sum_{r=1}^{n_l} \phi_{l,s,r}(z_{l,r}), \quad s = 1, \dots, n_{l+1}, \quad (3)$$

so that the whole network computes: $\text{KAN}(\mathbf{x}) = (\Phi_{L-1} \circ \Phi_{L-2} \circ \dots \circ \Phi_0)(\mathbf{x})$. Each $\phi_{l,s,r}$ is parameterized as a smooth B-spline,

$$\phi(z) = w_b b(z) + w_s \text{spline}(z),$$

with $b(z)$ a fixed baseline (e.g., identity or SiLU), and w_s, w_b are learnable weights. This construction is fully differentiable and trainable with standard backpropagation. While the shallow KART decomposition may involve irregular univariates, stacking layers yields smooth, accurate approximations (Liu et al., 2024b).

KANs retain universal approximation: refining spline grids (internal degrees of freedom) and stacking layers (external degrees of freedom) expands expressivity while exposing interpretable one-dimensional components. Compared to MLPs with fixed nonlinearities, complexity shifts from dense weight matrices to a small number of spline coefficients per edge, making the active graph (which edges matter) separable from the functional form of each transformation. Training follows standard optimizers and losses, with sparsity encouraged by ℓ_1 or group penalties that prune low-contribution edges. The pruned networks are compact, with remaining splines straightforward to inspect or approximate symbolically. The *MultKAN* extension (Liu et al., 2024a) further augments summation nodes with parameter-free multiplication nodes, explicitly representing interactions without forcing splines to emulate them.

In summary, KANs are deep models built from learnable one-dimensional splines composed through sums (and optionally products). They preserve universal approximation, train with off-the-shelf methods, admit effective pruning, and expose interpretable building blocks at the level of univariate functions and explicit interactions. These properties make them natural candidates to replace neural backbones in causalKANs.

4 CAUSAL KOLMOGOROV–ARNOLD NETWORKS

We introduce *causalKANs*, a framework that replaces the neural components of standard potential–outcome regressors with Kolmogorov–Arnold Networks (KANs), and augments training with *pruning* and *auto-symbolic search* to yield analytic, interpretable CATE formulas. The approach is model-agnostic: any causalNN can be converted by swapping MLP blocks with KAN (or MultKAN) blocks while keeping inputs and outputs unchanged, and using the same loss functions as their neural counterparts but augmented with regularization terms.

4.1 INTERPRETABLE CAUSAL LEARNING

Our goal is to deliver *closed-form* and *auditable* estimates of $\hat{\tau}(\mathbf{x}) = \hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$ without sacrificing the flexibility of deep learning. To this end, we expose a pipeline in which the practitioner controls each simplification step, chooses which steps to perform, and sets explicit performance budgets that cap the error introduced by simplification. The five stages are shown in Fig. 19; architecture substitution is illustrated in Fig. 1, and implementation details are in App. B.

- KAN-ification (architecture swap).** Choose a causalNN (e.g., S-KAN, T-KAN, TARKAN, DragonKAN) and replace each MLP subnetwork by a KAN block with matching input/output dimensions. When the original network has multiple subnetworks, intermediate widths and representation sizes need not be preserved; KAN widths can be reduced to encourage parsimony. MultKAN nodes are optional and only used when explicit interactions are needed (§3.2).
- Hyperparameter selection & training.** Tune depth, width, spline grid size, and regularization. During training we employ: **i)** *edge activity* penalty (ℓ_1) $\lambda_1 \sum_{l,s,r} \mathbb{E} |\phi_{l,s,r}(z_{l,r})|$ to promote sparse parent sets; **ii)** *spline coefficient* regularization $\lambda_c \sum |w_s| + \lambda_s \sum |w_s - w'_s|$ to shrink magnitudes and encourage smoothness; and **iii)** an *entropy* penalty on fan-in/fan-out to discourage diffuse connectivity. These terms stabilize pruning and symbolification. Early stopping is applied

on a validation split. For notational brevity, all penalties are aggregated into \mathcal{R} in later formulas. Among models with statistically indistinguishable validation metrics, we select the simplest (fewer layers, coarser grids, fewer nodes, no multiplication nodes), which is supported by an ablation study in App. C.2. We also compare the predictive loss (excluding \mathcal{R}) to the original causalNN; if the causalKAN loss exceeds the baseline by more than a user budget Λ_{arch} , we revise the KAN design.

3. **Pruning (structure simplification).** On a held-out set, compute edge importance scores

$$s_{l,s,r} = \mathbb{E} [|\phi_{l,s,r}(z_{l,r})|], \quad (4)$$

remove edges (and isolated nodes) with $s_{l,s,r} < \Gamma$, and retrain briefly if desired. There are alternative pruning techniques, that remove edges based on other metrics, see (Liu et al., 2024a). This step is optional: if the held-out loss rises beyond a pruning budget Λ_{prune} , we revert the change. Pruning exposes smaller parent sets and simplifies the subsequent symbolification.

4. **Auto-symbolic search (function simplification).** For each remaining edge function, fit a simple atom from a dictionary (identity, polynomials, tanh, sin, cos, log, exp, etc.) via

$$\min_{m,a,b,c,d} \frac{1}{|\mathcal{V}|} \sum_{z \in \mathcal{V}} \left(\phi_{l,s,r}(z) - [e f_m(az + b) + d] \right)^2. \quad (5)$$

We introduce an interpretability-first inductive bias: attempt the simplest atoms first (polynomials); if the R^2 exceeds a threshold Γ_{R^2} , accept immediately without testing more complex atoms. Otherwise, continue through the dictionary and accept a substitution only if the validation loss increases by at most Λ_{symb} . If the increase exceeds the budget, revert. This design makes the loss-simplicity trade-off explicit and user-controllable.

5. **CATE extraction and interpretation.** Compose the simplified univariate functions to obtain closed-form heads $\hat{\mu}_0(\mathbf{x})$ and $\hat{\mu}_1(\mathbf{x})$ and thus an explicit $\hat{\tau}(\mathbf{x})$, which we also simplify algebraically. These expressions are executable and auditable, and they can be inspected term by term or transported across settings by substituting domain-grounded functions if needed.

At stages (3) and (4) we enforce an accept-reject gate: a structural or symbolic change is kept only if the held-out performance stays within its budget ($\Lambda_{\text{prune}}, \Lambda_{\text{symb}}$); otherwise it is reverted. Practitioners may also choose to *skip* pruning and/or symbolification entirely, yielding a continuum from fully flexible KANs (budgets set to 0) to sparse, fully symbolified formulas (larger budgets). This is crucial in regulated or high-stakes settings: one can freeze the pipeline at the desired interpretability level and document any accuracy impact.

When $L = 1$ and no MultKAN nodes are used, each head reduces to a Kolmogorov-Arnold Additive Model (KAAM). In this setting, we provide **i) probability radar plots** (PRPs) summarizing each variable’s contribution relative to the average outcome, and **ii) partial dependence plots** (PDPs) showing the variation of the outcome as a function of a single covariate (Almodóvar et al., 2025; Knottenbelt et al., 2024). For deeper KANs or any use of MultKAN nodes, we currently provide only the closed-form expressions; visualization tools for these more complex models are left as future work. A detailed description and examples for KAAM-based plots appear in App. C.3.

In summary, the pipeline makes the accuracy-interpretability trade-off *controlled and reproducible*. Budgets ($\Lambda_{\text{arch}}, \Lambda_{\text{prune}}, \Lambda_{\text{symb}}$) bound the deviation from the original predictive performance; every accepted change is logged, and every rejected change is reverted. The procedure is architecture-agnostic, produces executable closed-form $\hat{\mu}_0(\mathbf{x})$, $\hat{\mu}_1(\mathbf{x})$, and $\hat{\tau}(\mathbf{x})$, and can be halted at any stage depending on the practitioner’s needs. We adopt this protocol across all causalKANs variants and experiments that follow.

Beyond the universal approximation, KANs are most useful when their inductive biases match the data generating process, in particular when response surfaces are smooth and approximately decomposable into low dimensional additive components. This is consistent with evidence from physics informed and PDE benchmarks, where KANs outperform MLPs under smooth, structured dynamics (Wang et al., 2025), and with the original KAN work, which argues that many real world systems admit sparse and smooth functional structure (Liu et al., 2024b). From a causal perspective, treatment effect functions are often simpler than the underlying potential outcomes (Curth and Van der Schaar, 2021a), so causalKANs can leverage pruning and symbolic substitution as implicit regularizers that bias

toward simple, interpretable effect formulas, in contrast to standard MLP backbones. This behavior is illustrated in our SCM experiment (Section 5.3), where symbolic KANs recover the correct smooth additive structure and outperform both MLPs and non symbolic KANs, and suggests that causalKANs are particularly suitable in domains where smoothness and approximate additivity are plausible, such as physics based models, biomedical dose–response, or structured policy applications.

4.2 CAUSALKANS VARIANTS

We instantiate four canonical architectures, standard baselines in the CATE literature, to enable fair comparisons with prior work. These are representative examples: the same substitution process applies to other causalNNs, see the App. B.1 for an extended discussion. Fig. 3 shows a schematic view where each block can be a KAN, and App. B details specific implementations.

S-KAN (*causalKAN for S-Learner*) uses a single KAN to predict outcomes:

$$\hat{\mu}(\mathbf{x}, t) = \text{KAN}(\mathbf{x}, t), \quad \hat{\mu}_0(\mathbf{x}) = \hat{\mu}(\mathbf{x}, 0), \quad \hat{\mu}_1(\mathbf{x}) = \hat{\mu}(\mathbf{x}, 1). \quad (6)$$

With one KAN layer, Eq. 6 reduces to an additive model (KAAM). S-KAN naturally supports continuous treatments, since $\mu_t(\mathbf{x}, \mathbf{t})$ can be evaluated for $\mathbf{t} \in \mathbb{R}$.

T-KAN (*causalKAN for T-Learner*) employs two independent KAN heads:

$$\hat{\mu}_0(\mathbf{x}) = \text{KAN}_0(\mathbf{x}), \quad \hat{\mu}_1(\mathbf{x}) = \text{KAN}_1(\mathbf{x}). \quad (7)$$

Each head is updated only with its respective treatment group, using

$$\mathcal{L}_{\text{T-KAN}} = \frac{1}{N} \sum_{i=1}^N \left[t_i (y_i - \hat{\mu}_1(\mathbf{x}_i))^2 + (1 - t_i) (y_i - \hat{\mu}_0(\mathbf{x}_i))^2 \right] + \mathcal{R}. \quad (8)$$

This setup allows treatment-specific fits but may increase variance due to data splitting. Either head can be restricted to KAAM for maximal simplicity.

TARKAN (*causalKAN for TARNet*) first computes a representation vector,

$$\mathbf{z}(\mathbf{x}) = \text{KAN}_z(\mathbf{x}) \in \mathbb{R}^{D_z}, \quad (9)$$

then feeds it to two KAN heads as in T-KAN :

$$\hat{\mu}_0(\mathbf{x}) = \text{KAN}_0(\mathbf{z}(\mathbf{x})), \quad \hat{\mu}_1(\mathbf{x}) = \text{KAN}_1(\mathbf{z}(\mathbf{x})). \quad (10)$$

Training follows the T-style loss Eq. 8, applied after the representation. This design parallels TARNet and relates to Mehendale et al. (2025), though we do not impose explicit distribution-matching on \mathbf{z} .

DragonKAN (*causalKAN for DragonNet*) augments TARKAN with a propensity head:

$$\hat{e}(\mathbf{x}) = \sigma(\text{KAN}_e(\mathbf{z}(\mathbf{x}))), \quad (11)$$

where σ is the sigmoid function. Training adds a cross-entropy penalty on $\hat{e}(\mathbf{x})$ to the outcome losses in Eq. 8, encouraging \mathbf{z} to capture confounding. We omit targeted regularization layers to preserve formula simplicity, focusing on potential–outcome regression and CATE estimation. Unlike S-KAN, these architectures (T-KAN, TARKAN, DragonKAN) are limited to binary or discrete treatments.

5 EXPERIMENTAL EVALUATION

We evaluate *causalKANs* on semi-synthetic benchmarks where ground-truth potential outcomes are available: ACIC 2016 (Dorie et al., 2019) and IHDP (Hill, 2011), settings A and B. Semi-synthetic data allow objective assessment of CATE accuracy despite the fundamental problem of causal inference. We report **i**) the Precision in Estimation of Heterogeneous Effect (PEHE) and **ii**) ATE error.

Given test set $\mathcal{D}_{\text{test}}$, PEHE is $\sqrt{\frac{1}{|\mathcal{D}_{\text{test}}|} \sum (\hat{\tau}(x_i) - \tau(x_i))^2}$, and ATE error is $|\hat{\tau} - \tau|$, where $(\hat{\tau}, \tau)$ are the estimated and the real ATE, computed as $\tau = \frac{1}{|\mathcal{D}_{\text{test}}|} \sum \tau(x)$. All training details follow §4.

Datasets. IHDP (A/B). The IHDP benchmark (Hill, 2011) is constructed from covariates of a real observational study, with synthetic outcomes generated from known functions. We use the standard

100 replications for settings A and B, where A corresponds to a linear outcome surface and B introduces nonlinear and heterogeneous effects. *ACIC 2016*. The ACIC’16 challenge datasets (Dorie et al., 2019) provide covariates from administrative health records and semi-synthetic treatments and outcomes generated through nonlinear mechanisms. We use the nonlinear regimes commonly adopted in prior work. *NSLM*. The NSLM benchmark uses covariates from the National Study of Learning Mindset (Yeager et al., 2019), with a semi-synthetic DGP introduced by Carvalho et al. (2019). *NEWS*. The NEWS dataset contains text-derived covariates from a corpus of 5000 documents (Johansson et al., 2016). We adopt the DGP proposed by Crabbé et al. (2022). *TCGA*. The TCGA benchmark uses RNA-seq covariates from the Cancer Genome Atlas (Weinstein et al., 2013), combined with the semi-synthetic DGP of Zhang et al. (2023). Further details for all datasets are provided in App. C.1.

5.1 COMPARISON WITH NEURAL NETS

Baselines. We benchmark S-/T-KAN, TARKAN, and DragonKAN against their MLP-based counterparts (S-/T-learner, TarNet, DragonNet). We intentionally restrict baselines to neural counterparts to test our central claim—*comparable accuracy with improved interpretability*—rather than to chase leaderboards.

We trained causalNNs with the same hyperparameter budget as causalKANs, modifying depth, number of neurons, activations, regularization and learning rate, and selected the best model based on validation loss (we justify this choice in App. C.2).

Tab 1 reports PEHE and ATE error per dataset/setting. We assess statistical significance using the Friedman test with Holm-corrected Wilcoxon signed-rank post-hoc comparisons at $\alpha = 0.05$ (see Demšar (2006); Rainio et al. (2024) for a details on these tests). In the tables, the best-performing model (baseline) is underlined. Models in **bold** are those for which we cannot reject the null hypothesis of equal performance with respect to the baseline ($p \geq 0.05$). A train/val/test split of 80/10/10 was leveraged, as well as early stopping and Adam optimizer (Kingma and Ba, 2014).

We observe from both Tab 1 and Fig. 2 that, in 7 of the 8 datasets, there exists at least one instance of causalKAN that is the best model or whose performance is statistically indistinguishable from the best-performing neural model, in the sense that we fail to reject the null hypothesis that the models achieve the same value of the evaluation metric. On the other hand, in the TCGA dataset, none of the causalKANs achieve competitive results. However, that does not prevent the use of our pipeline, since one of the steps of it is to compare the metrics against causalNNs. If metrics of causalKANs are not satisfactory, then the use of KAN based predictors is not recommended by our pipeline (see line 14 of Alg. 1).

5.2 INTERPRETABILITY RESULTS

We show here some examples that highlight the interpretability of causalKANs, and we will refer to App. C.3 for more examples and details on these representations. We selected randomly one realization of the shown datasets.

Table 1: **Out-of-sample** ATE error and PEHE for KAN and MLP across datasets. The baseline (best) according to the Friedman corrected test is underlined, and all models not statistically different are **bolded** ($\alpha \geq 0.05$ in a paired Wilcoxon corrected test). Values are reported as mean_{std} .

Dataset	Model	KAN		MLP	
		ATE err	PEHE	ATE err	PEHE
IHDP A	S-Learner	0.19 _{0.35}	0.98 _{1.03}	0.23 _{0.47}	1.04 _{1.30}
	T-Learner	0.13 _{0.08}	0.62 _{0.28}	0.24 _{0.27}	0.90 _{0.51}
	TarNet	0.13 _{0.10}	0.64 _{0.31}	0.17 _{0.27}	0.70 _{0.78}
	DragonNet	0.14 _{0.10}	0.66 _{0.35}	0.17 _{0.27}	0.68 _{0.77}
IHDP B	S-Learner	0.37 _{0.37}	3.01 _{0.63}	0.32 _{0.24}	2.63 _{0.58}
	T-Learner	0.34 _{0.27}	2.80 _{0.44}	0.25 _{0.20}	2.06 _{0.35}
	TarNet	0.33 _{0.24}	2.68 _{0.44}	0.23 _{0.20}	2.08 _{0.36}
	DragonNet	0.28 _{0.22}	2.64 _{0.44}	0.25 _{0.21}	2.00 _{0.34}
ACIC 2	S-Learner	0.20 _{0.38}	0.20 _{0.38}	0.38 _{0.44}	0.74 _{0.38}
	T-Learner	0.56 _{0.42}	1.43 _{0.72}	1.97 _{2.13}	7.05 _{6.48}
	TarNet	0.25 _{0.33}	0.78 _{0.44}	0.36 _{0.44}	0.96 _{0.83}
	DragonNet	0.21 _{0.26}	0.75 _{0.33}	0.36 _{0.44}	0.96 _{0.82}
ACIC 7	S-Learner	0.66 _{0.75}	4.13 _{1.56}	0.75 _{0.60}	4.51 _{1.54}
	T-Learner	0.43 _{0.43}	3.06 _{1.18}	1.37 _{2.14}	7.05 _{6.64}
	TarNet	0.42 _{0.43}	3.06 _{1.18}	0.58 _{0.50}	4.17 _{1.47}
	DragonNet	0.43 _{0.43}	3.06 _{1.17}	0.59 _{0.50}	4.17 _{1.48}
ACIC 26	S-Learner	0.42 _{0.45}	3.23 _{1.57}	0.75 _{0.60}	4.51 _{1.54}
	T-Learner	0.36 _{0.34}	2.80 _{1.08}	1.37 _{2.14}	7.05 _{6.64}
	TarNet	0.35 _{0.34}	2.80 _{1.08}	0.58 _{0.50}	4.17 _{1.47}
	DragonNet	0.35 _{0.34}	2.79 _{1.08}	0.59 _{0.50}	4.17 _{1.48}
NSLM	S-Learner	0.048 _{0.038}	0.752 _{0.008}	0.190 _{0.032}	0.753 _{0.006}
	T-Learner	0.050 _{0.033}	0.752 _{0.006}	0.048 _{0.033}	0.756 _{0.007}
	TarNet	0.055 _{0.034}	0.752 _{0.006}	0.049 _{0.033}	0.755 _{0.007}
	DragonNet	0.058 _{0.040}	0.753 _{0.007}	0.050 _{0.035}	0.756 _{0.007}
NEWS	S-Learner	2.01 _{0.68}	0.141 _{0.119}	2.73 _{1.04}	0.284 _{0.311}
	T-Learner	2.04 _{0.53}	0.192 _{0.136}	2.12 _{0.54}	0.158 _{0.126}
	TarNet	1.97 _{0.48}	0.157 _{0.125}	2.05 _{0.44}	0.120 _{0.112}
	DragonNet	2.04 _{0.72}	0.176 _{0.125}	2.11 _{0.44}	0.128 _{0.112}
TCGA $\times 10^{-2}$	S-Learner	2.50 _{1.50}	4.64 _{1.26}	0.03 _{0.02}	2.18 _{0.27}
	T-Learner	1.98 _{6.36}	4.41 _{6.27}	0.02 _{0.01}	2.15 _{6.07}
	TarNet	1.19 _{0.92}	3.26 _{0.65}	0.04 _{0.01}	1.81 _{0.03}
	DragonNet	14.95 _{33.92}	20.34 _{33.81}	0.05 _{0.04}	3.74 _{0.27}

432 **ACIC-7.** We show results for a single layer
 433 T-KAN (as it is an additive model, we call it
 434 T-KAAM) for ACIC-7, which has been demon-
 435 strated very high-performance according to
 436 Tab 1. This model yields a function of *hetero-*
 437 *geneous* CATE, which is a generalized additive
 438 model of the covariates (see App. B). In addition
 439 to the closed-form CATE, we provide in Fig. 4
 440 (extended in App. C.3.1) a *Radar plot* (Left) that
 441 compares for two individuals, what is the contribu-
 442 tion of each variable to the CATE, compared
 443 with the average outcome, and a *PDP* (Right) that
 444 represent the curve that defines the variation of the
 445 CATE (Δ CATE) with a specific feature. In this case, we can observe how x_{44} increases the CATE
 w.r.t. the average in individual 10, while it decreases the CATE in individual 11, because the shape of
 the curve that relates x_{44} and CATE has a maximum near the value of x_{44} of individual 10.

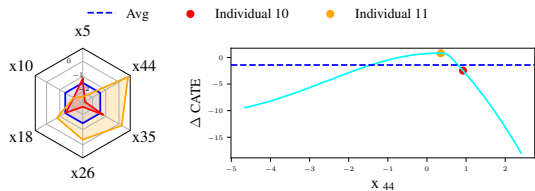


Figure 4: Radar plot and PDP for variable contribution to CATE, using T-KAAM in ACIC-7.

446 **IHDP A.** On the other hand, we can estimate *homogen-*
 447 *eous* CATEs with a shallow S-KAN (S-KAAM), which
 448 provides closed-form homogeneous treatment effect. In
 449 IHDP A, where the known (denoised) ITE is 4 for every
 450 individual, we leverage S-KAAM to obtain the value of
 451 the causal effect. In Fig. 5 we can observe that the pre-
 452 dicted formula has an additive linear term relative to the
 453 treatment, which correspond to the CATE (as developed in
 454 App. B). The contribution of other variables (not causal)
 455 can be consulted in App. C.3.1. In this case, the CATE can
 456 directly be consulted in the formula: 3.74 for the given
 457 example.

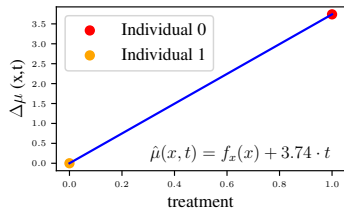


Figure 5: PDP for treatment contribution in $\hat{\mu}(x, t)$ estimation, using S-KAAM in IHDP-A.

5.3 SYMBOLIC CATE RECOVERY

460 We have conducted two experiments with known structural causal models, to determine empirically
 461 how well causalKANs capture the equations of the data generating process (DGP). They also provide
 462 examples of when the regularization and the simplification steps of the pipeline of causalKANs help
 463 to achieve better approximations to the true data generating functions. Metrics can be consulted in
 464 App. C.4

465 In Fig. 6a, the function that generates the potential outcomes is additive in the treatment. Therefore, it
 466 lies in the class that S-KAAM can model (see App. B). It can be observed that **i**) the S-KAAM provides
 467 a function that is closer to the groundtruth compared with MLP and **ii**) the symbolic substitution step
 468 approximates the groundtruth even better. Note that there is a constant value that S-KAAM cannot
 469 recover, but does not modify the CATE estimation, since each difference $\hat{f}(x, t_1) - \hat{f}(x, t_1)$ will
 470 cancel this constant.
 471

472 On the other hand, Fig. 6b represents the approximation to the CATE with a binary treatment, as the
 473 difference between the two potential outcomes, $\{y(1), y(0)\}$. In this case, individual treatment effect
 474 is not constant across individuals, and the treatment does not contribute additively to the outcomes.
 475 However, the contribution to each variable to the ITE is additive. Therefore, this DGP lies in the class
 476 that T-KAAM can model (see App. B). In the same fashion as in the previous experiment, symbolic
 477 T-KAAM is the model that captures better the true function that generates the CATE.
 478

6 FINAL REMARKS

481 We have presented causalKANs, a framework that transforms high-performing neural CATE esti-
 482 mators into knowledge-augmented networks with closed-form, auditable effect functions. On
 483 standard benchmarks, at least one causalKAN variant matches or surpasses its neural counterpart on
 484 PEHE/ATE, while shallow heads (KAAMs) often yield the most favorable accuracy-interpretability
 485 trade-off. The pruning and auto-symbolification stages expose explicit formulas and partial depend-
 ence plots (PDP; Friedman (2001)). These properties make causalKANs suitable for healthcare,

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

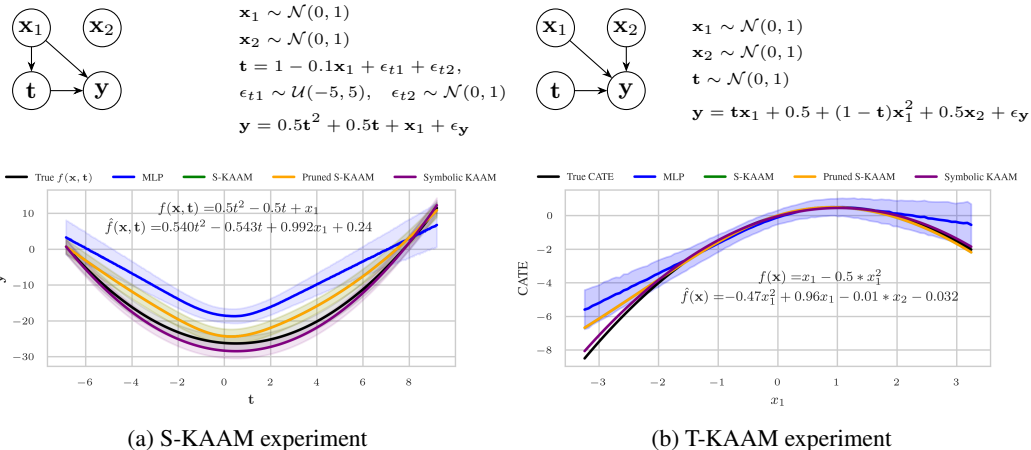


Figure 6: Synthetic SCMs and the curves recovered by causalKANs in each step of the pipeline. In blue, the curve recovered by MLP counterparts. Curves green and yellow are overlapped, because the effect of pruning is not appreciable. f and \hat{f} represent the groundtruth and the estimated function by Symbolic KAAM, respectively. All curves has been shifted so that the first point coincides. (Left) represents directly the equation of the outcome and covariates provided by S-KAAM. (Right) represents the equation of the CATE given by T-KAAM, as the difference of two additive formulas.

policy evaluation, and economics, where interpretability is central to adoption and auditing (Amann et al., 2020; Doshi-Velez and Kim, 2017; Rudin, 2018; Tonekaboni et al., 2019). Closed-form expressions can be scrutinized, simplified, or aligned with domain-grounded terms to support external validity. We release code to ensure reproducibility and facilitate adoption.

Nevertheless, causalKANs also present open challenges. Visualization tools are currently most effective for shallow KAAMs, while scaling them to deeper or more complex networks remains an area for development. Training and inference are generally more demanding than for standard neural models, and performance can be sensitive to regularization choices or pruning/symbolification budgets, which may introduce approximation errors or variability. As with all observational CATE estimation, our conclusions depend on standard identification assumptions (e.g., unconfoundedness, overlap), and our evaluation relies mainly on semi-synthetic datasets, so further assessment on field data is needed. Finally, uncertainty quantification for symbolic outputs and broader support for continuous or multi-valued treatments are still limited.

Building on these observations, future work should expand theoretical and empirical aspects of causalKANs. Methodologically, interpretability should be assessed with task-grounded and user-study metrics (Tonekaboni et al., 2019), and contrasted with post-hoc explanations such as LIME and SHAP (Lundberg and Lee, 2017; Ribeiro et al., 2016, routinely used but with different guarantees). Developing principled surrogate metrics for model selection, and procedures that reduce accuracy gaps across the simplification stages, would improve stability and reliability. Extending causalKANs to continuous and multi-treatment settings, as well as deriving finite-sample guarantees for CATE estimation, remain important directions. Incorporating inverse-probability weighting, doubly robust objectives, and batch sampling techniques would further broaden scope. Empirically, validation on real-world data, ideally against interventional. Another avenue is to adapt causalKANs to non-tabular modalities such as images, graphs, and genomic or sequence data by leveraging recent KAN variants for convolutional and graph architectures (Bodner and coauthors, 2024; Bresson and coauthors, 2024; Cherednichenko and coauthors, 2025; Liu et al., 2024b). Designing benchmarks with ground-truth potential outcomes in these modalities is itself an open challenge, and we leave multimodal CATE benchmarking with causalKANs to future work. Related with visualization, beyond PDPs, implementing related plot families, such as accumulated local effects (Apley and Zhu, 2020) or individual conditional expectation (Goldstein et al., 2015), may also be informative when applied directly to the closed-form formulas. Finally, exploring robustness to adversarial perturbations, distribution shifts, and fairness constraints (Rudin, 2018), together with reducing the performance gaps between KAN-ification, pruning, and symbolification, will be important for reliable deployment.

REPRODUCIBILITY STATEMENT

We have taken several steps to ensure the reproducibility of our results. All causalKAN variants are defined with explicit mathematical formulations in §4.2, and the interpretability pipeline is described in detail in §4.1, including pruning and auto-symbolic search criteria, budgets, and regularization terms. The assumptions and problem setup are clearly stated in §3, while implementation details such as hyperparameter choices, training protocols, and evaluation procedures are provided in §5 and the appendix. All experiments were run on publicly available benchmark datasets. To facilitate reproducibility and comparison with future work, we release the full source code, including scripts to reproduce all figures and tables. Together, these resources ensure that both the methodological contributions and empirical findings of this paper can be independently verified and extended by the community.

ETHICS STATEMENT

This work does not raise direct ethical concerns of that nature. All datasets used are publicly available benchmarks, processed following documented procedures. The proposed framework aims to improve transparency and interpretability in causal machine learning, which is a key requirement for responsible deployment of AI systems in high-stakes domains (Amann et al., 2020; Rudin, 2018). By providing closed-form analytic formulas and explicit control over the trade-off between accuracy and interpretability, our approach is designed to facilitate auditing, external validation, and trustworthy adoption. As with all machine learning methods, potential biases or limitations in the datasets may propagate to the models, but the interpretability of our approach makes such issues easier to identify and mitigate. This work adheres to the ICLR code of ethics.

BIBLIOGRAPHY

- R. Agarwal, L. Melnick, N. Frosst, X. Zhang, B. Lengerich, R. Caruana, and G. E. Hinton. Neural additive models: Interpretable machine learning with neural nets. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 4699–4711. Curran Associates, Inc., 2021. (Cited in page 2.)
- A. Al Bataineh, B. Vamsi, M. El-Abd, and B. P. Doppala. Kolmogorov-arnold networks for predicting carotid intima-media thickness in cardiovascular risk assessment. *Scientific Reports*, 15(1):32108, 2025. (Cited in page 2.)
- A. Alaa and M. Schaar. Limits of estimating heterogeneous treatment effects: Guidelines for practical algorithm design. In *International Conference on Machine Learning*, pages 129–138. PMLR, 2018. (Cited in page 4.)
- F. Alharbi, N. Budhiraja, A. Vakanski, B. Zhang, M. Elbashir, H. Guduru, and M. Mohammed. Interpretable graph kolmogorov-arnold networks for multi-cancer classification and biomarker identification using multi-omics data. *Scientific Reports*, 15, 07 2025. doi: 10.1038/s41598-025-13337-0. (Cited in page 2.)
- A. Almodóvar, P. A. Apellániz, A. Garrido, F. Fernández-Salvador, S. Zazo, and J. Parras. Interpretable clinical classification with kolgomorov-arnold networks. *arXiv preprint arXiv:2509.16750*, 2025. (Cited in pages 2, 6, 19, and 25.)
- J. Amann, A. Blasimme, E. Vayena, D. Frey, and V. I. Madai. Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC Medical Informatics and Decision Making*, 20(310), 2020. doi: 10.1186/s12911-020-01332-6. (Cited in pages 1, 10, and 11.)
- D. W. Apley and J. Zhu. Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(4):1059–1086, 2020. doi: 10.1111/rssb.12377. (Cited in page 10.)
- V. I. Arnold. On functions of three variables. *Doklady Akademii Nauk SSSR*, 114:679–681, 1957. English translation in *American Mathematical Society Translations*, Series 2, Vol. 28: Sixteen Papers on Analysis, 1963, pp. 51–54. (Cited in page 4.)

- 594 S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the*
595 *National Academy of Sciences*, 113(27):7353–7360, 2016. doi: 10.1073/pnas.1510489113. (Cited
596 in page 2.)
- 597 H. Bang and J. M. Robins. Doubly robust estimation in missing data and causal inference models.
598 *Biometrics*, 61(4):962–973, 2005. (Cited in page 21.)
- 600 F. J. Bargagli-Stoffi, R. Cadei, K. Lee, and F. Dominici. Causal rule ensemble: Interpretable discovery
601 and inference of heterogeneous treatment effects. *arXiv preprint arXiv:2009.09036*, 2020. (Cited
602 in page 3.)
- 603 A. D. Bodner and coauthors. Convolutional kolmogorov–arnold networks. *arXiv preprint*
604 *arXiv:2406.13155*, 2024. (Cited in page 10.)
- 606 A. D. Bodner, A. S. Tepsich, J. N. Spolski, and S. Pourteau. Convolutional kolmogorov-arnold
607 networks. *arXiv preprint arXiv:2406.13155*, 2024. (Cited in page 2.)
- 608 J. Braun and M. Griebel. On a constructive proof of kolmogorov’s superposition theorem. *Construct-*
609 *ive approximation*, 30(3):653–675, 2009. (Cited in page 4.)
- 611 R. Bresson and coauthors. Kolmogorov–arnold networks meet graph learning. *Transactions on*
612 *Machine Learning Research*, 2024. (Cited in page 10.)
- 613 R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. Intelligible models for healthcare:
614 Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM*
615 *SIGKDD international conference on knowledge discovery and data mining*, pages 1721–1730,
616 2015. (Cited in page 2.)
- 617 C. Carvalho, A. Feller, J. Murray, S. Woody, and D. Yeager. Assessing treatment effect variation in
618 observational studies: Results from a data challenge. *Observational Studies*, 5(2):21–35, 2019.
619 (Cited in pages 2, 8, and 22.)
- 621 K. Chen, Q. Yin, and Q. Long. Covariate-balancing-aware interpretable deep learning models for
622 treatment effect estimation. *Statistics in Biosciences*, 17(1):132–150, April 2025. doi: 10.1007/s1
623 2561-023-09394-6. (Cited in page 3.)
- 624 O. Cherednichenko and coauthors. Kolmogorov–arnold networks for genomic tasks. *Briefings in*
625 *Bioinformatics*, 2025. (Cited in page 10.)
- 627 S.-Y. Cho, S. Lee, and H.-G. Kim. Forecasting vix using interpretable kolmogorov-arnold networks.
628 *Expert Systems with Applications*, 294:128781, 2025. ISSN 0957-4174. doi: <https://doi.org/10.1016/j.eswa.2025.128781>. (Cited in page 2.)
- 629 J. Crabbé, A. Curth, I. Bica, and M. van der Schaar. Benchmarking heterogeneous treatment effect
630 models through the lens of interpretability, 2022. (Cited in pages 8 and 22.)
- 631 A. Curth and M. Van der Schaar. On inductive biases for heterogeneous treatment effect estimation.
632 *Advances in Neural Information Processing Systems*, 34:15883–15894, 2021a. (Cited in pages 4
633 and 6.)
- 634 A. Curth and M. Van der Schaar. Nonparametric estimation of heterogeneous treatment effects: From
635 theory to learning algorithms. In *International Conference on Artificial Intelligence and Statistics*,
636 pages 1810–1818. PMLR, 2021b. (Cited in pages 4 and 24.)
- 637 A. Curth, R. W. Peck, E. McKinney, J. Weatherall, and M. van der Schaar. Using machine learning to
638 individualize treatment effect estimation: Challenges and opportunities. *Clinical Pharmacology &*
639 *Therapeutics*, 115(4):710–719, 2024. doi: 10.1002/cpt.3159. (Cited in page 1.)
- 640 J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning*
641 *research*, 7(Jan):1–30, 2006. (Cited in page 8.)
- 642 V. Dorie, J. Hill, U. Shalit, M. Scott, and D. Cervone. Automated versus do-it-yourself methods for
643 causal inference: Lessons learned from a data analysis competition. 2019. (Cited in pages 2, 7, 8,
644 and 22.)

- 648 F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. *arXiv*
649 *preprint arXiv:1702.08608*, 2017. (Cited in pages 2 and 10.)
650
- 651 European Parliament and Council. Regulation (EU) 2016/679 (GDPR) on the protection of natural
652 persons with regard to the processing of personal data and on the free movement of such data.
653 Official Journal of the European Union, 2016. (Cited in page 1.)
654
- 655 European Parliament and Council. Regulation (EU) 2024/1689 (Artificial Intelligence Act). Official
656 Journal of the European Union, 2024. Transparency and information requirements for high-risk AI
657 systems. (Cited in page 1.)
- 658 J. C. Foster, J. M. Taylor, and S. J. Ruberg. Subgroup identification from randomized clinical trial
659 data. *Statistics in Medicine*, 30, 2011. (Cited in page 2.)
660
- 661 D. Frauen and S. Feuerriegel. Estimating individual treatment effects under unobserved confounding
662 using binary instruments. *arXiv preprint arXiv:2208.08544*, 2022. (Cited in page 4.)
663
- 664 J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*,
665 pages 1189–1232, 2001. (Cited in pages 9 and 25.)
666
- 667 A. Goldstein, A. Kapelner, J. Bleich, and E. Pitkin. Peeking inside the black box: Visualizing
668 statistical learning with plots of individual conditional expectation. *journal of Computational and*
Graphical Statistics, 24(1):44–65, 2015. (Cited in pages 10 and 25.)
669
- 670 B. Goodman and S. Flaxman. European union regulations on algorithmic decision-making and a
671 “right to explanation”. *AI magazine*, 38(3):50–57, 2017. (Cited in page 1.)
672
- 673 P. Hahn, C. Carvalho, D. Puelz, and J. He. Regularization and confounding in linear regression
674 for treatment effect estimation. *Bayesian Analysis*, 13(1):163–182, 2018. ISSN 1936-0975. doi:
675 10.1214/16-BA1044. Publisher Copyright: © 2018 International Society for Bayesian Analysis.
(Cited in page 2.)
676
- 677 M. A. Hernán and J. M. Robins. *Causal Inference: What If*. CRC Press, Boca Raton, 1st edition,
678 2025. ISBN 978-0367711337. (Cited in pages 3 and 20.)
679
- 680 J. L. Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and*
681 *Graphical Statistics*, 20(1):217–240, 2011. doi: 10.1198/jcgs.2010.08162. (Cited in pages 2, 7,
and 22.)
682
- 683 K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approxi-
684 mators. *Neural networks*, 2(5):359–366, 1989. (Cited in page 4.)
685
- 686 K. Imai and A. Strauss. Estimation of heterogeneous treatment effects from randomized experiments,
687 with application to the optimal planning of the get-out-the-vote campaign. *Political Analysis*, 19
(1):1–19, 2011. doi: 10.1093/pan/mpq035. (Cited in page 1.)
688
- 689 F. D. Johansson, U. Shalit, and D. Sontag. Learning representations for counterfactual inference.
690 In *Proceedings of the 33rd International Conference on International Conference on Machine*
691 *Learning - Volume 48*, ICML’16, page 3020–3029. JMLR.org, 2016. (Cited in pages 2, 8, and 22.)
692
- 693 E. H. Kennedy. Towards optimal doubly robust estimation of heterogeneous causal effects. *Electronic*
694 *Journal of Statistics*, 14(1):38–56, 2020. (Cited in page 21.)
695
- 696 D. M. Kent, E. Steyerberg, and D. van Klaveren. Personalized evidence based medicine: predictive
697 approaches to heterogeneous treatment effects. *BMJ*, 363:k4245, 2018. doi: 10.1136/bmj.k4245.
(Cited in page 1.)
698
- 699 C. Kim and O. Bastani. Learning interpretable models with causal guarantees. *arXiv preprint*
700 *arXiv:1901.08576*, 2019. (Cited in page 3.)
701
- 702 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
(Cited in page 8.)

- 702 W. Knottenbelt, Z. Gao, R. Wray, W. Z. Zhang, J. Liu, and M. Crispin-Ortuzar. Coxkan:
703 Kolmogorov-arnold networks for interpretable, high-performance survival analysis. *arXiv preprint*
704 *arXiv:2409.04290*, 2024. (Cited in page 6.)
705
- 706 A. N. Kolmogorov. On the representation of continuous functions of several variables as superposi-
707 tions of continuous functions of a smaller number of variables. *Doklady Akademii Nauk SSSR*, 108
708 (2):179–182, 1956. (Cited in page 4.)
- 709 A. N. Kolmogorov. On the representations of continuous functions of many variables by superposition
710 of continuous functions of one variable and addition. In *Dokl. Akad. Nauk USSR*, volume 114,
711 pages 953–956, 1957. (Cited in page 4.)
712
- 713 S. R. Künzel, J. S. Sekhon, P. J. Bickel, and B. Yu. Metalearners for estimating heterogeneous
714 treatment effects using machine learning. *Proceedings of the national academy of sciences*, 116
715 (10):4156–4165, 2019. (Cited in pages 1, 4, and 21.)
- 716 Y. Ling, P. Upadhyaya, L. Chen, X. Jiang, and Y. Kim. Emulate randomized clinical trials using
717 heterogeneous treatment effect estimation for personalized treatments: Methodology review and
718 benchmark. *Journal of biomedical informatics*, 137:104256, 2023. (Cited in page 4.)
719
- 720 Z. Liu, P. Ma, Y. Wang, W. Matusik, and M. Tegmark. Kan 2.0: Kolmogorov-arnold networks meet
721 science. *arXiv preprint arXiv:2408.10205*, 2024a. (Cited in pages 1, 4, 5, and 6.)
- 722 Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljačić, T. Y. Hou, and M. Tegmark. Kan:
723 Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756*, 2024b. (Cited in pages 1, 2, 3, 4,
724 5, 6, 10, 26, and 30.)
725
- 726 J. Loftus, L. Bynum, and S. Hansen. Causal dependence plots. *Advances in Neural Information*
727 *Processing Systems*, 37:112656–112683, 2024. (Cited in page 25.)
- 728 S. Lolak, J. Attia, G. J. McKay, and A. Thakkinstian. Application of dragonnet and conformal
729 inference for estimating individualized treatment effects for personalized stroke prevention: Retro-
730 spective cohort study. *JMIR cardio*, 9:e50627, 2025. (Cited in page 4.)
731
- 732 S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in*
733 *Neural Information Processing Systems*, pages 4765–4774, 2017. (Cited in page 10.)
- 734 C. F. Manski. Statistical treatment rules for heterogeneous populations. *Econometrica*, 72(4):
735 1221–1246, 2004. (Cited in page 1.)
736
- 737 E. Mehendale, A. Thorat, R. Kolla, and N. Pedanekar. Kanite: Kolmogorov-arnold networks for ite
738 estimation. *arXiv preprint arXiv:2503.13912*, 2025. (Cited in pages 3 and 7.)
- 739 M. Moradi, S. Panahi, E. Bollt, and Y.-C. Lai. Kolmogorov-arnold network autoencoders. *arXiv*
740 *preprint arXiv:2410.02077*, 2024. (Cited in page 2.)
741
- 742 R. Moraffah, M. Karami, R. Guo, A. Raglin, and H. Liu. Causal interpretability for machine learning
743 - problems, methods and evaluation. *SIGKDD Explor. Newsl.*, 22(1):18–33, May 2020. ISSN
744 1931-0145. doi: 10.1145/3400051.3400058. (Cited in page 3.)
- 745 X. Nie and S. Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*, 108
746 (2):299–319, 2021. (Cited in pages 4 and 21.)
747
- 748 O. H. M. Padilla, Y. Chen, and G. Ruiz. A causal fused lasso for interpretable heterogeneous treatment
749 effects estimation. *arXiv preprint arXiv:2110.00901*, 2021. (Cited in page 3.)
- 750 J. Pearl. *Causality*. Cambridge university press, 2009. (Cited in page 25.)
751
- 752 V. S. Pendyala and N. Venkatachalam. The effectiveness of kolmogorov–arnold networks in the
753 healthcare domain. *Applied Sciences*, 15(16), 2025. ISSN 2076-3417. (Cited in page 2.)
754
- 755 O. Rainio, J. Teuvo, and R. Klén. Evaluation metrics and statistical tests for machine learning.
Scientific Reports, 14(1):6086, 2024. (Cited in page 8.)

- 756 M. T. Ribeiro, S. Singh, and C. Guestrin. “Why Should I Trust You?” explaining the predictions of
757 any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge*
758 *Discovery and Data Mining*, pages 1135–1144, 2016. doi: 10.1145/2939672.2939778. (Cited in
759 page 10.)
- 760 J. M. Robins, A. Rotnitzky, and L. P. Zhao. Estimation of regression coefficients when some regressors
761 are not always observed. *Journal of the American Statistical Association*, 89(427):846–866, 1994.
762 (Cited in page 21.)
- 763 P. R. Rosenbaum and D. B. Rubin. The central role of the propensity score in observational studies
764 for causal effects. *Biometrika*, 70(1):41–55, 1983. (Cited in page 20.)
- 765 D. B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies.
766 *Journal of educational Psychology*, 66(5):688, 1974. (Cited in page 3.)
- 767 C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use
768 interpretable models instead. *Nature Machine Intelligence*, 1:206 – 215, 2018. (Cited in pages 3,
769 10, and 11.)
- 770 M. J. Saary. Radar plots: a useful way for presenting multivariate health care data. *Journal of clinical*
771 *epidemiology*, 61(4):311–317, 2008. (Cited in page 24.)
- 772 P. Sanchez, J. P. Voisey, T. Xia, H. I. Watson, A. Q. O’Neil, and S. A. Tsaftaris. Causal machine
773 learning for healthcare and precision medicine. *Royal Society Open Science*, 9(8):220638, 2022.
774 (Cited in page 1.)
- 775 M. Schmidt and H. Lipson. Distilling free-form natural laws from experimental data. *Science*, 324
776 (5923):81–85, 2009. doi: 10.1126/science.1165893. (Cited in page 2.)
- 777 P. Schwab, L. Linhardt, and W. Karlen. Perfect match: A simple method for learning representations
778 for counterfactual inference with neural networks. *arXiv preprint arXiv:1810.00656*, 2018. (Cited
779 in pages 2 and 4.)
- 780 P. Schwab, L. Linhardt, S. Bauer, J. M. Buhmann, and W. Karlen. Learning counterfactual represent-
781 ations for estimating individual dose-response curves. In *Proceedings of the AAAI Conference on*
782 *Artificial Intelligence*, volume 34, pages 5612–5619, 2020. (Cited in page 2.)
- 783 U. Shalit, F. D. Johansson, and D. Sontag. Estimating individual treatment effect: generalization
784 bounds and algorithms. In *International conference on machine learning*, pages 3076–3085.
785 PMLR, 2017. (Cited in pages 1, 2, 3, 4, and 24.)
- 786 C. Shi, D. Blei, and V. Veitch. Adapting neural networks for the estimation of treatment effects.
787 *Advances in neural information processing systems*, 32, 2019. (Cited in pages 1, 4, and 24.)
- 788 H. Shuai and F. Li. Physics-informed kolmogorov-arnold networks for power system dynamics.
789 *IEEE Open Access Journal of Power and Energy*, 12:46–58, 2025. ISSN 2687-7910. doi:
790 10.1109/oajpe.2025.3529928. (Cited in page 2.)
- 791 J. Splawa-Neyman, D. M. Dabrowska, and T. P. Speed. On the application of probability theory to
792 agricultural experiments. essay on principles. section 9. *Statistical Science*, pages 465–472, 1990.
793 (Cited in page 3.)
- 794 G. Tesei, S. Giampanis, J. Shi, and B. Norgeot. Learning end-to-end patient representations through
795 self-supervised covariate balancing for causal treatment effect estimation. *Journal of biomedical*
796 *informatics*, 140:104339, 2023. (Cited in page 4.)
- 797 S. Tonekaboni, S. Joshi, M. D. McCradden, and A. Goldenberg. What clinicians want: Contextualiz-
798 ing explainable machine learning for clinical end use. In *Proceedings of the 4th Machine Learning*
799 *for Healthcare Conference*, volume 106, pages 359–380. PMLR, 2019. (Cited in pages 1 and 10.)
- 800 A. A. Tsiatis. *Semiparametric Theory and Missing Data*. Springer, New York, 2006. (Cited in
801 page 21.)

- 810 M. J. Van der Laan, S. Rose, et al. *Targeted learning: causal inference for observational and*
811 *experimental data*, volume 4. Springer, 2011. (Cited in page 4.)
812
- 813 S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using random
814 forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018. doi: 10.1080/
815 01621459.2017.1319839. (Cited in page 1.)
- 816 H. Wakaura, R. Mulyawan, and A. B. Suksmono. Enhanced variational quantum kolmogorov-arnold
817 network. *arXiv preprint arXiv:2503.22604*, 2025. (Cited in page 2.)
818
- 819 Y. Wang, J. Sun, J. Bai, C. Anitescu, M. S. Eshaghi, X. Zhuang, T. Rabczuk, and Y. Liu. Kolmogorov-
820 arnold-informed neural network: A physics-informed deep learning framework for solving forward
821 and inverse problems based on kolmogorov-arnold networks. *Computer Methods in Applied*
822 *Mechanics and Engineering*, 433:117518, 2025. (Cited in pages 3 and 6.)
- 823 J. N. Weinstein, E. A. Collisson, G. B. Mills, K. R. Shaw, B. A. Ozenberger, K. Ellrott, I. Shmulevich,
824 C. Sander, and J. M. Stuart. The cancer genome atlas pan-cancer analysis project. *Nature genetics*,
825 45(10):1113–1120, 2013. (Cited in pages 8 and 22.)
- 826 Y. Wu, H. Liu, K. Ren, and X. Chang. Causal rule learning: Enhancing the understanding of
827 heterogeneous treatment effect via weighted causal rules. *arXiv preprint arXiv:2310.06746*, 2023.
828 (Cited in page 3.)
829
- 830 L. Yao, S. Li, Y. Li, M. Huai, J. Gao, and A. Zhang. Representation learning for treatment effect
831 estimation from observational data. *Advances in neural information processing systems*, 31, 2018.
832 (Cited in pages 2 and 4.)
- 833 D. S. Yeager, P. Hanselman, G. M. Walton, J. S. Murray, R. Crosnoe, C. Muller, E. Tipton,
834 B. Schneider, C. S. Hulleman, C. P. Hinojosa, et al. A national experiment reveals where a
835 growth mindset improves achievement. *Nature*, 573(7774):364–369, 2019. (Cited in pages 8 and
836 22.)
837
- 838 J. Yoon, J. Jordon, and M. van der Schaar. GANITE: Estimation of individualized treatment effects
839 using generative adversarial nets. In *International Conference on Learning Representations*, 2018.
840 (Cited in page 2.)
- 841 R. C. Yu, S. Wu, and J. Gui. Residual kolmogorov-arnold network for enhanced deep learning. *arXiv*
842 *preprint arXiv:2410.05500*, 2024. (Cited in page 2.)
843
- 844 R. Zhang, A. Lensen, and C. Shang. Speeding up genetic programming based symbolic regression
845 using gpus. In *Proc. of EvoApplications*, 2022. (Cited in page 2.)
- 846 Y. Zhang, H. Zhang, Z. C. Lipton, L. E. Li, and E. Xing. Exploring transformer backbones for
847 heterogeneous treatment effect estimation. *Transactions on Machine Learning Research*, 2023.
848 ISSN 2835-8856. (Cited in pages 3, 8, and 22.)
- 849 Q. Zhao and T. Hastie. Causal interpretations of black-box models. *Journal of Business & Economic*
850 *Statistics*, 39(1):272–281, 2021. (Cited in page 25.)
851
852
853
854
855
856
857
858
859
860
861
862
863

Appendix

A TREATMENT SPACE GENERALIZATIONS

Let us generalize the theory explained in §3 to the settings of multiple treatment and continuous treatment. CausalKANs can also be generalized following the same fashion.

A.1 MULTIPLE DISCRETE TREATMENTS

We extend the setup in §3 to $\mathbf{t} \in \mathcal{T} = \{1, \dots, K\}$. For each $k \in \{1, \dots, K\}$, define the potential outcome $\mathbf{y}(k)$ and the potential outcome regressor

$$\mu_k(\mathbf{x}) := \mathbb{E}[\mathbf{y}(k) \mid \mathbf{x} = \mathbf{x}]. \quad (12)$$

Pairwise conditional average treatment effects are contrasts

$$\tau_{ba}(\mathbf{x}) := \mu_b(\mathbf{x}) - \mu_a(\mathbf{x}), \quad a, b \in \{1, \dots, K\}. \quad (13)$$

Assumptions generalize in the standard way: **i**) positivity, $P(\mathbf{t} = k \mid \mathbf{x} = \mathbf{x}) > 0$ for all k ; **ii**) conditional ignorability, $\mathbf{y}(k) \perp\!\!\!\perp \mathbf{t} \mid \mathbf{x}$ for all k ; **iii**) consistency, $\mathbf{y}_i(t_i) = y_i$; and **iv**) no interference. Let the propensity vector be $e_k(\mathbf{x}) = P(\mathbf{t} = k \mid \mathbf{x} = \mathbf{x})$, $\sum_k e_k(\mathbf{x}) = 1$.

S-learner (unchanged). Use a single regressor $\hat{\mu}(\mathbf{x}, t)$; predicted heads are $\hat{\mu}_k(\mathbf{x}) = \hat{\mu}(\mathbf{x}, k)$. Train with factual MSE:

$$\mathcal{L}_{\text{S-Learner}} = \frac{1}{N} \sum_{i=1}^N (\hat{\mu}(\mathbf{x}_i, t_i) - y_i)^2. \quad (14)$$

T-learner (K heads). Instantiate K heads $\hat{\mu}_k(\mathbf{x})$ and update only the factual head:

$$\mathcal{L}_{\text{T-learner}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu}_{t_i}(\mathbf{x}_i))^2. \quad (15)$$

TARNet (K heads). Use a shared representation $\mathbf{z}(\mathbf{x}) \in \mathbb{R}^{D_z}$ and K heads:

$$\mathbf{z}(\mathbf{x}), \quad \hat{\mu}_k(\mathbf{z}(\mathbf{x})) = \text{Head}_k(\mathbf{z}(\mathbf{x})), \quad \mathcal{L}_{\text{TARNet}} = \frac{1}{N} \sum_i (y_i - \hat{\mu}_{t_i}(\mathbf{z}(\mathbf{x}_i)))^2. \quad (16)$$

DragonNet ($K+1$ heads). Add a propensity head with softmax output:

$$\hat{e}_k(\mathbf{z}(\mathbf{x})) = \text{softmax}_k(\mathbf{z}(\mathbf{x})), \quad \mathcal{L}_{\text{DragonNet}} = \mathcal{L}_{\text{TARNet}} + \lambda_{PS} \left[-\frac{1}{N} \sum_i \log \hat{e}_{t_i}(\mathbf{z}(\mathbf{x}_i)) \right]. \quad (17)$$

After training, pairwise CATEs follow from equation 13 with μ_k replaced by $\hat{\mu}_k$.

A.2 CONTINUOUS TREATMENTS

Let $\mathbf{t} \in \mathcal{T} \subset \mathbb{R}$ and define the dose–response function

$$\mu(\mathbf{x}, t) := \mathbb{E}[\mathbf{y}(t) \mid \mathbf{x} = \mathbf{x}]. \quad (18)$$

For any reference level $t_0 \in \mathcal{T}$, the conditional effect of moving from t_0 to t is

$$\tau(\mathbf{x}, t) := \mu(\mathbf{x}, t) - \mu(\mathbf{x}, t_0), \quad (19)$$

and local effects can be summarized by the marginal treatment effect $\partial\mu(\mathbf{x}, t)/\partial t$ if desired. The identifiability assumptions extend as: positivity w.r.t. the treatment density $p(t \mid \mathbf{x})$, conditional ignorability $\mathbf{y}(t) \perp\!\!\!\perp t \mid \mathbf{x}$ for all $t \in \mathcal{T}$, consistency, and no interference.

Applicable architecture. Among the causalNNs above, only the S-learner directly accommodates continuous t :

$$\mathcal{L}_{\text{S-cont}} = \frac{1}{N} \sum_{i=1}^N (\hat{\mu}(\mathbf{x}_i, t_i) - y_i)^2, \tag{20}$$

with $\hat{\tau}(\mathbf{x}, t) = \hat{\mu}(\mathbf{x}, t) - \hat{\mu}(\mathbf{x}, t_0)$. $\hat{\tau}(\mathbf{x}, t)$ can also be provided directly as a function of the treatment. T-learner/TARNet/DragonNet require finitely many treatment-specific heads and hence are not directly applicable without discretization of \mathcal{T} ; we therefore recommend the use of the S-learner in continuous-treatment settings.

B DETAILS ON CAUSALKANS VARIANTS

As explained in §4, the adaptation of causalNNs to causalKANs consists of changing the MLP backbones by KANs. In Fig. 7, we include examples of the proposed causalKANs, with an arbitrary number of layers and summation nodes.

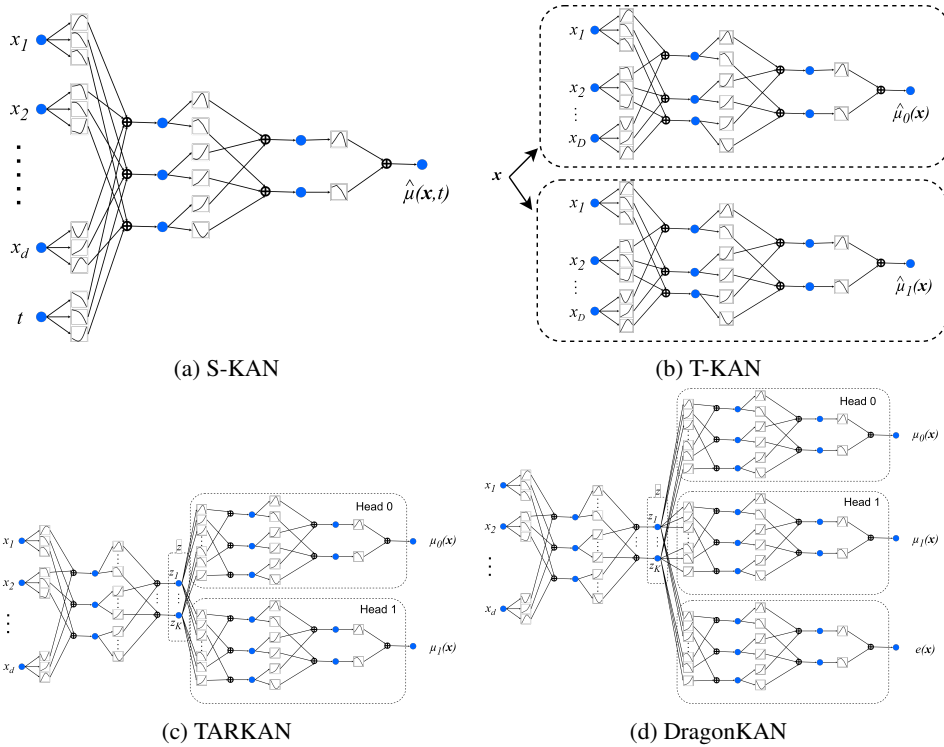


Figure 7: CausalKANs detailed architectures. The number of layers of each model is arbitrary, just to show an example of each one.

From the point of view of interpretability, constructing very complex causalKANs is harmful for achieving simple and auditable formulas. Therefore, as mentioned in point 2 of §4.1, obtain the simplest model is crucial to improve interpretability, and we recommend to minimize the complexity of the network, among the models with similar performance.

We want to illustrate some interesting combinations of hyperparameters, that yield into special interpretable causal effects. In particular, when meta-learners employ additive models, we call them S-KAAM and T-KAAM, and when the heads of TARKAN and DragonKAN are also shallow additive models, we call them TARKAAM and DragonKAAM.

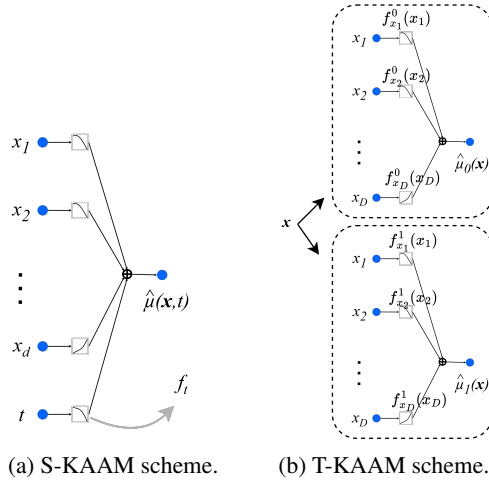


Figure 8: Particularly interpretable architectures of causalKANs. (Left) a single KAAM that predicts the outcome in an additive manner. (Left) two KAAMs, one for each treatment, which leads to an additive CATE function.

S-KAAM. When training and S-KAN, if the hyperparameters selected results into a single-layer KAN, then we have a KAAM (Almodóvar et al., 2025). We call this model S-KAAM, and a scheme can be found in Fig. 8a. The interesting point is that, when we have an S-KAAM, the effects on the population are homogeneous, and we have access to the *effect curve* directly. The effect curve defines how the outcome varies with the treatment. Since, in a S-KAAM, the model is additive:

$$\hat{\mu}(\mathbf{x}, t) = f_{\mathbf{x}}(\mathbf{x}) + f_t(t), \quad (21)$$

we define the *effect curve*, $f_t : \mathcal{T} \rightarrow \mathbb{R}$, as the function that represents the variations of the outcome, depending on the treatment. Therefore, any causal effect of two different values of the treatment, a, b can be computed as the difference along the effect curve:

$$\hat{\tau}_{ab}(\mathbf{x}) = \hat{\mu}(\mathbf{x}, b) - \hat{\mu}(\mathbf{x}, a) = [f_{\mathbf{x}}(\mathbf{x}) + f_t(b)] - [f_{\mathbf{x}}(\mathbf{x}) + f_t(a)] = f_t(b) - f_t(a), \quad (22)$$

Note that this model yields homogeneous treatment effects, making the CATE equivalent to the ATE, and removing the covariate dependence.

For a binary treatment, the CATE can be computed as:

$$\hat{\tau}(\mathbf{x}) = f_t(1) - f_t(0), \quad (23)$$

T-KAAM. In the same fashion, if both subnetworks in a T-KAN are KAAMs, we can get a simple analytic solution of the causal effect, that depends on each covariate independently.

Having that each potential outcome can be expressed as:

$$\hat{\mu}_t(\mathbf{x}) = f_{x_1}^t(x_1) + f_{x_2}^t(x_2) + \dots + f_{x_D}^t(x_D), \quad (24)$$

where $f_{x_i}^t$ is the activation function corresponding to the covariate x_i in the subnetwork of the treatment t . We can compute the individual contribution of each variable, yielding in an estimation of the causal effect, for a binary treatment, as:

$$\hat{\tau}(\mathbf{x}) = \underbrace{[f_{x_1}^1(x_1) - f_{x_1}^0(x_1)]}_{\mathbf{g}_{x_1}(x_1)} + \underbrace{[f_{x_2}^1(x_2) - f_{x_2}^0(x_2)]}_{\mathbf{g}_{x_2}(x_2)} + \dots + \underbrace{[f_{x_D}^1(x_D) - f_{x_D}^0(x_D)]}_{\mathbf{g}_{x_D}(x_D)} \quad (25)$$

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

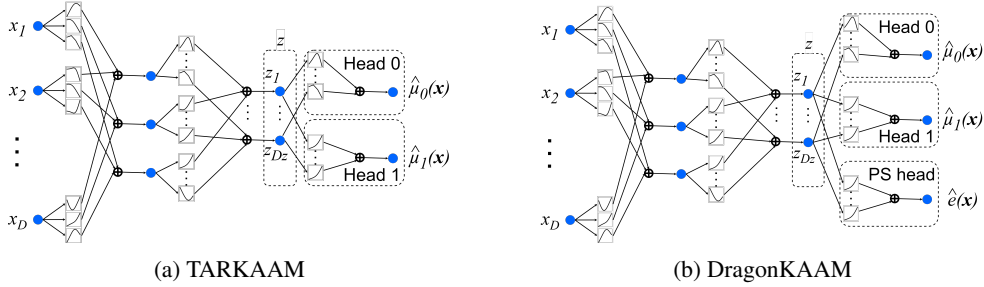


Figure 9: TARKAN and DragonKAN with shallow heads. With respect to Fig. 7, the connections that depart from \mathbf{z} have been reordered to make the equivalence more noticeable.

As an implementation detail, note that a T-KAAM can be implemented with a single KAN shallow model with two outputs, each one for $\hat{\mu}_1(\mathbf{x})$ and $\hat{\mu}_0(\mathbf{x})$, respectively.

TARKAAM and DragonKAAM. We would like to make an observation on the implementation of TARKAN and DragonKAN when the heads are shallow KAAMs. We can observe in Fig. 9 that, when instantiating a single KAN, with 2 and 3 output nodes, respectively. To construct those figures, we have reordered the last layer of a single KAN with 2/3 outputs, and show that is equivalent to add 2/3 KAAMs that have $\mathbf{z}(\mathbf{x})$ as input.

That is an interesting fact because we find, empirically, that TARKAAM and DragonKAAM are between the best-performer TAR-like and Dragon-like networks when evaluating in the IHDP and ACIC datasets. We also note that a KAN can learn identity splines, which could lead to have a representation vector equal to the input: $\mathbf{z}(\mathbf{x}) = \mathbf{x}$. If that is the case, the TARKAAM and DragonKAAM can be seen as T-KAAMs, since all the outputs are independent additive functions of the covariates (except for the common regularization term, \mathcal{R}). For example, that is the case for the dataset ACIC-7, which cause that T-KAN, TARKAN and DragonKAN to have the almost same metrics.

B.1 EXTENDING CAUSALKANS TO OTHER ESTIMATORS

CausalKANs can be used as a drop in replacement in many causal estimators by KAN-ifying nuisance components (outcome regressions, propensity scores, CATE regressors) while leaving the estimand unchanged. After pruning and symbolification, these components become analytic expressions that can be inspected and manipulated. Which part should be made interpretable is estimator dependent and also dataset dependent, since different applications emphasize different aspects of the causal mechanism.

We consider binary treatment $\mathbf{t} \in \{0, 1\}$, outcome \mathbf{y} , covariates \mathbf{x} , propensity score $e(\mathbf{x}) = \mathbb{P}(\mathbf{t} = 1 \mid \mathbf{x})$, potential outcome regressions $\mu_0(\mathbf{x})$ and $\mu_1(\mathbf{x})$, and CATE $\tau(\mathbf{x}) = \mu_1(\mathbf{x}) - \mu_0(\mathbf{x})$ with estimator $\hat{\tau}(\mathbf{x})$.

IPW. Inverse probability weighting (IPW) estimates the ATE using only a propensity model $e(\mathbf{x})$ (Hernán and Robins, 2025; Rosenbaum and Rubin, 1983). For a sample $\{(\mathbf{x}_i, t_i, y_i)\}_{i=1}^n$,

$$\hat{\tau}_{\text{IPW}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{t_i y_i}{\hat{e}(\mathbf{x}_i)} - \frac{(1 - t_i) y_i}{1 - \hat{e}(\mathbf{x}_i)} \right). \quad (26)$$

causalKANs enter through

$$\hat{e}(\mathbf{x}) = \sigma(f_{\text{KAN}}^{(e)}(\mathbf{x})), \quad (27)$$

so that the log odds

$$\log \frac{\hat{e}(\mathbf{x})}{1 - \hat{e}(\mathbf{x})} = f_{\text{KAN}}^{(e)}(\mathbf{x}) \quad (28)$$

is a closed form KAN expression. After pruning and symbolification, this gives an interpretable model of the selection mechanism (which covariates drive treatment assignment and how they interact), while $\hat{\tau}_{\text{IPW}}$ remains the same functional of the weights.

X learner. The X learner (Künzel et al., 2019) combines outcome regressions and effect regressions. Given $\hat{\mu}_0(\mathbf{x})$ and $\hat{\mu}_1(\mathbf{x})$, one constructs pseudo effects for treated and control units and learns two CATE regressors $\hat{\tau}(\mathbf{x})^{(1)}$ and $\hat{\tau}(\mathbf{x})^{(0)}$, which are combined as

$$\hat{\tau}(\mathbf{x})_{\text{X}} = g(\mathbf{x}) \hat{\tau}(\mathbf{x})^{(0)} + (1 - g(\mathbf{x})) \hat{\tau}(\mathbf{x})^{(1)}, \quad (29)$$

with $g(\mathbf{x})$ often chosen as $\hat{e}(\mathbf{x})$.

CausalKANs can be used in two complementary ways:

- *KAN-ified potential outcomes:* set $\hat{\mu}_t(\mathbf{x}) = f_{\text{KAN}}^{(t)}(\mathbf{x})$ for $t \in \{0, 1\}$. After pruning and symbolification, $\hat{\mu}_0(\mathbf{x})$ and $\hat{\mu}_1(\mathbf{x})$ are analytic functions, so the practitioner can directly inspect how covariates shape each potential outcome and the induced pseudo effects.
- *KAN-ified CATE:* set $\hat{\tau}(\mathbf{x})^{(t)} = g_{\text{KAN}}^{(t)}(\mathbf{x})$ and optionally $g(\mathbf{x}) = \hat{e}(\mathbf{x}) = \sigma(f_{\text{KAN}}^{(e)}(\mathbf{x}))$. Then $\hat{\tau}(\mathbf{x})_{\text{X}}$ is an explicit combination of a small number of KAN terms, providing a closed form heterogeneous effect curve.

Which variant is preferable depends on whether interpretability should focus on the potential outcomes, the CATE, or the selection mechanism, and this is driven by the dataset and scientific question.

AIPW, DR learners, and R learner. Augmented IPW (AIPW) combines outcome and propensity models and is doubly robust (Bang and Robins, 2005; Robins et al., 1994; Tsiatis, 2006). A standard AIPW ATE estimator is

$$\hat{\tau}_{\text{AIPW}} = \frac{1}{n} \sum_{i=1}^n \left[\hat{\mu}_1(\mathbf{x}_i) - \hat{\mu}_0(\mathbf{x}_i) + \frac{t_i(y_i - \hat{\mu}_1(\mathbf{x}_i))}{\hat{e}(\mathbf{x}_i)} - \frac{(1 - t_i)(y_i - \hat{\mu}_0(\mathbf{x}_i))}{1 - \hat{e}(\mathbf{x}_i)} \right]. \quad (30)$$

KAN-ifying $\hat{\mu}_0$, $\hat{\mu}_1$ and/or \hat{e} yields interpretable models for baseline heterogeneity $\hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$ and for the augmentation terms. For applications where the effect structure is primary, causalKANs on $\hat{\mu}_0$ and $\hat{\mu}_1$ are most informative; for settings where selection bias is central, causalKANs on \hat{e} are more relevant.

DR learners and R-type learners construct pseudo outcomes and then regress them on \mathbf{x} using a flexible CATE model (e.g. Kennedy, 2020; Nie and Wager, 2021). In these methods, causalKANs are naturally used as the final CATE regressor, that is

$$\hat{\tau}(\mathbf{x}) = g_{\text{KAN}}^{(\tau)}(\mathbf{x}), \quad (31)$$

trained with an orthogonalized loss. After pruning and symbolification, the resulting $\hat{\tau}(\mathbf{x})$ is a closed form expression for the heterogeneous effect that inherits the robustness properties of the underlying DR or R-type learner. If desired, causalKANs can also parameterize nuisance components such as \hat{e} or $\hat{\mu}_t$, but this is optional and should be guided by which functions the practitioner wishes to interpret.

Summary. These examples illustrate a general pattern: causalKANs can replace the neural building blocks of many estimators (IPW, meta learners, DR and R learners, etc.) while leaving the estimand (ATE or CATE) unchanged. After pruning and symbolification, the replaced components become analytic, auditable functions of \mathbf{x} . The most relevant use of causalKANs is estimator dependent (propensity vs potential outcomes vs CATE) and also dataset dependent, since different applications require understanding different parts of the causal pipeline.

1134 C DETAILED RESULTS

1135

1136

1137

C.1 DETAILS OF THE DATASETS

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

We have conducted a systematic process of experiments on several datasets. We include here some details of the datasets, and their procedence, for completeness.

IHDP. (Hill, 2011) The Infant Health and Development Program (IHDP) is a semi-synthetic benchmark derived from a real observational study on the effect of early childhood intervention. The publicly used version was constructed by Hill (2011), who retained 25 pre-treatment covariates and simulated potential outcomes to enable pointwise ground truth for individual treatment effects. Following standard practice, we use the canonical 100 replications that define settings A and B. Setting A yields homogeneous treatment effects via linear outcome surfaces, while setting B introduces nonlinear relationships and heterogeneous effects. The dataset contains 747 units with a binary treatment and substantial treatment imbalance, since all treated units from one site were removed in the construction. Covariates include demographic, prenatal, and birth-related variables. All outcomes in the benchmark are synthetic and do not correspond to clinical endpoints.

ACIC 2016. (Dorie et al., 2019) The ACIC 2016 causal inference challenge provides a large collection of semi-synthetic datasets derived from real covariates in the Louisiana Medicaid program. Each instance consists of approximately $N \approx 4800$ units and 58 covariates, including demographic, socioeconomic, and healthcare-related variables. Treatments and potential outcomes are generated through complex nonlinear mechanisms with covariate-dependent assignment, enabling rigorous benchmarking of estimation procedures under realistic confounding. We focus on the nonlinear polynomial and exponential outcome regimes (for example, settings 2, 7, and 26), which are commonly used in prior work. For each selected setting, we use the 77 replications released as part of the competition, each providing a distinct draw of treatment assignment and potential outcomes while keeping covariates fixed across replications. As in other semi-synthetic benchmarks, the outcomes are simulated and have no policy or medical interpretation.

NSLM. The dataset contains covariables from the national study of learning mindset (Yeager et al., 2019). We used the DGP followed by Carvalho et al. (2019), where there are 10000 datapoints (10000 students across 76 schools), with 3 categorical variables and 6 continuous variables, and the treatment and the outcome are simulated with no meaning.

TCGA. The Cancer Genomic Atlas (TCGA) (Weinstein et al., 2013). The TCGA project collected gene expression data from various types of cancers in 9659 individuals. In this case, we conducted the data generating process proposed by Zhang et al. (2023), which keeps only the data from 100 covariates of the RNA sequence. The outcome is simulated and does not have physical meaning. Although in the paper of Zhang et al. (2023), they explore several treatments and dosage, we restrict only a binary treatment without dosage, to have a fair comparison between all the models evaluated.

NEWS. The NEWS dataset (Johansson et al., 2016) was created to perform causal inference. The covariates come from a text database, of 5000 documents, and they represent the count of each word. In total, there are 3477 possible words in each datapoint. The treatment is the device in which the users read, and the outcome is simulated and represent the reader experience. However, we followed the DGP from Crabbé et al. (2022), that employs 100 components of the principal component analysis of the covariates to generate the potential outcomes.

C.2 ABLATION STUDIES

Complexity. We further analyze the relationship between model complexity and the precision in estimation of heterogeneous effects (PEHE). To this end, we define a *complexity score* that aggregates the contributions of different architectural and regularization choices. Specifically, hidden dimensions contribute 0 if no hidden layers are used, 2 if a single hidden layer of size 5 is used, and 3 otherwise. The weight penalty λ contributes 0 if set to 0.01 and 1 if set to 0.001, since more regularization yields simpler models. The spline grid contributes 1, 2, or 3 for grid sizes $\{1, 3, 5\}$, respectively. Similarly, the polynomial order k contributes 1, 2, or 3 for $k \in \{1, 3, 5\}$, and sparse initialization reduces the score by one unit. Hence, the complexity score captures the combined effect of the number of hidden units, regularization strength, spline resolution, polynomial order, and initialization strategy. In addition, we compute the number of hidden layers directly from the hidden dimensions: 0 layers (no hidden units, corresponding to KAAM), 1 layer (single hidden layer), and 2 layers (deeper networks).

We then represent the correlation between both complexity and number of hidden layers against PEHE, using linear regression fits and Pearson correlation coefficients across all hyperparameter configurations (Figs. 10 and 11). The results indicate that the correlation is weak: the fitted slopes are close to zero and the Pearson coefficients are generally small. Interestingly, increasing the number of hidden layers often leads to slightly higher PEHE values, although the effect is minor and not consistent across all datasets. Similarly, increasing the overall complexity score does not systematically reduce PEHE.

These findings suggest that, in general, simple KAN architectures achieve competitive PEHE performance, and that increasing architectural or regularization complexity does not yield clear improvements. However, note that this score is a purely descriptive simple heuristic designed to summarize hyperparameter choices and also prioritize more interpretable models.

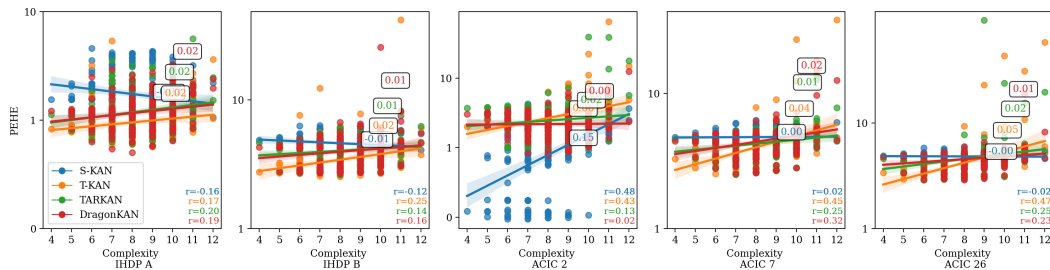


Figure 10: Correlation between model complexity and PEHE across datasets. Each point corresponds to a trained model with a given complexity score. Regression fits (with slopes annotated) and Pearson correlation coefficients are shown. The weak correlations indicate that increasing complexity does not improve PEHE.

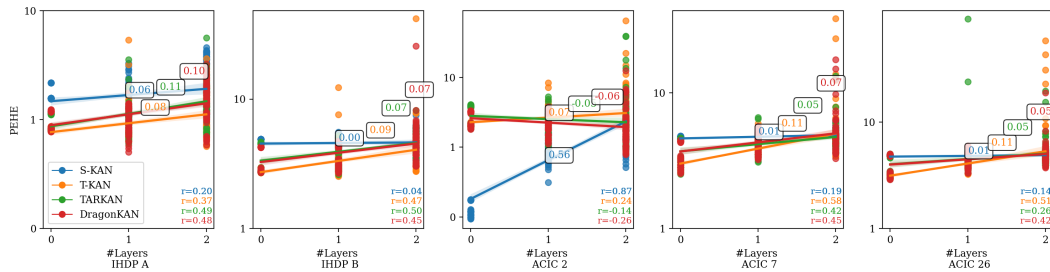
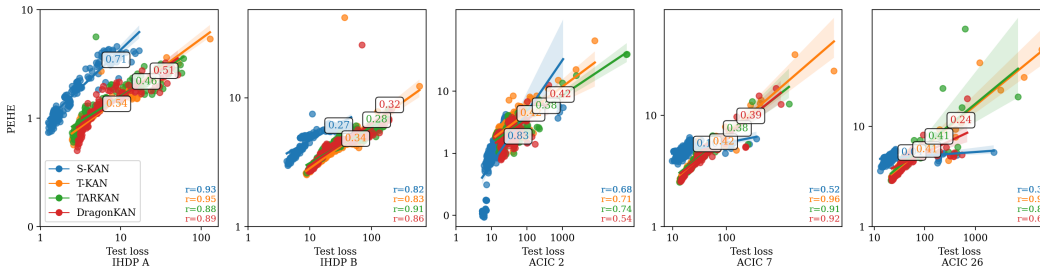


Figure 11: Correlation between the number of hidden layers and PEHE across datasets. The number of layers is computed from the hidden dimensions: 0 (no hidden layers, i.e., KAAM), 1 (single hidden layer), and 2 (two hidden layers). Regression fits and Pearson correlations show that deeper models tend to slightly increase PEHE, although the effect is minor.

1242 **Hyperparameter selection with Loss.** The hyperparameter selection is made based on the val-
 1243 idation loss. We want to show that the predictive test loss (i.e. the loss of the KAN model without
 1244 including sparsity regularizations) is a good proxy for selecting hyperparameters, since the PEHE
 1245 cannot be computed during evaluation because the real ITE is not known in real-world data and can
 1246 only be computed with (semi)synthetic data. We can observe in Fig. 12 that these quantities are
 1247 linearly correlated. However, the study of other surrogate metrics for model selection is still a topic
 1248 for future work.



1259 Figure 12: Regression plots of PEHE and Test predictive loss (logarithmic scale) for each dataset
 1260 with all causalKANs. Slope (in boxes), Pearson coefficient, r , (in bottom right corners) and 95% CI
 1261 (translucent bands) reported. We can observe a clear linear correlation between PEHE and Test loss
 1262 for all models and all datasets.

1264 Despite this correlation, we acknowledge that outcome-prediction loss is not, in general, a perfect
 1265 proxy for CATE accuracy, but all models evaluated in our work are potential-outcome regressors
 1266 trained solely on factual losses. For this family of estimators, the standard practice is to select
 1267 hyperparameters using validation prediction loss, since these methods do not expose orthogonalized
 1268 or doubly robust objectives (Shalit et al., 2017; Shi et al., 2019). To ensure a fair comparison,
 1269 causalKANs use the same training and selection rules.

1270 In Fig. 12, the plateau reflects intrinsic properties of outcome-regression methods, not an issue
 1271 introduced by causalKANs. Still, we observe strong correlations (typically $r > 0.8$) between factual
 1272 test loss and PEHE, indicating that loss remains a model-selection proxy for this class of learners. The
 1273 slight flattening for the models is consistent with known behavior of non-orthogonalized estimators
 1274 (Curth and Van der Schaar, 2021b; Shalit et al., 2017), where outcome prediction can improve faster
 1275 than CATE error.

1277 C.3 INTERPRETABILITY VISUALIZATIONS

1279 First of all, let us explain in detail the different visualization tools, in addition to the formula analysis
 1280 *per se*, that will help us to understand outcome variations, depending of the model used.

1282 **Probability Radar Plots (PRPs).** Probability Radar Plots (PRPs) (Saary, 2008) provide an inter-
 1283 pretable visualization of generalized additive models (GAMs) by mapping the isolated contribution
 1284 of each covariate into a radial plot. In our setting, each component function depends on a single
 1285 covariate, which enables a decomposition of the prediction into additive terms. We denote this
 1286 decomposition by Δ .

$$\Delta = \begin{bmatrix} f_{x_{1,1}} & f_{x_{1,2}} & \cdots & f_{x_{1,D}} \\ f_{x_{2,1}} & f_{x_{2,2}} & \cdots & f_{x_{2,D}} \\ \vdots & \vdots & \ddots & \vdots \\ f_{x_{N,1}} & f_{x_{N,2}} & \cdots & f_{x_{N,D}} \end{bmatrix}, \tag{32}$$

1293 where $f_{x_{i,j}}$ denotes the contribution of the feature j for individual i .

1294 To construct a PRP, we first compute the average contribution of each covariate across all individuals,
 1295 providing a baseline that reflects the average outcome (for S-KAAM) or the average conditional

average treatment effect (for T-KAAM). Then, for a given individual i , we plot the vector of deviations

$$\left(f_{x_{i,1}} - \frac{1}{N} \sum_{\ell=1}^N f_{x_{\ell,1}}, \dots, f_{x_{i,D}} - \frac{1}{N} \sum_{\ell=1}^N f_{x_{\ell,D}}\right), \quad (33)$$

which highlights how the contribution of each covariate for individual i differs from the population average. By arranging these deviations radially, PRPs enable intuitive comparison across covariates and between individuals.

Partial dependence plots (PDPs). We employ partial dependence plots (PDPs) (Friedman, 2001) as visualization tools for interpreting causalKANs. Unlike their conventional use in a “black-box” fashion—where the model is queried without access to its internals—KAN-based PDPs (Almodóvar et al., 2025) directly exploit the analytic equations of the learned model, displaying the underlying splines that constitute the predictors. This provides more faithful and transparent representations of the learned dependencies.

As noted by Loftus et al. (2024), PDPs that vary one covariate in isolation may be misleading in settings with mediators, since they ignore induced changes in other covariates. In our case, however, the assumptions in §3 guarantee that the covariates form a valid adjustment set, excluding mediators. Under these conditions, PDPs are valid tools to represent causal dependencies (Zhao and Hastie, 2021), and in fact rely on the same backdoor adjustment formula (Pearl, 2009) when the effect is homogeneous.

Formally, PDPs display the average variation in the predicted outcome as one covariate is varied, marginalizing over the remaining covariates. Their individual-level counterpart, individual conditional expectation (ICE) curves (Goldstein et al., 2015), provide counterfactual explanations for specific samples. Under additivity, PDPs and ICE coincide; under non-additivity, PDPs correspond to averages of the heterogeneous ICE curves.

C.3.1 HETEROGENEOUS CATE WITH T-KAAM

We show an example of how T-KAAM captures the CATE as a closed formula that can be interpreted. For the dataset *ACIC 7*, the T-KAAM model is one of the best-performers, and we can observe that the CATE can be represented as an addition of functions of each variable independently (see Eq. 24). An example obtained with one realization of the dataset can be observed in the following equation.

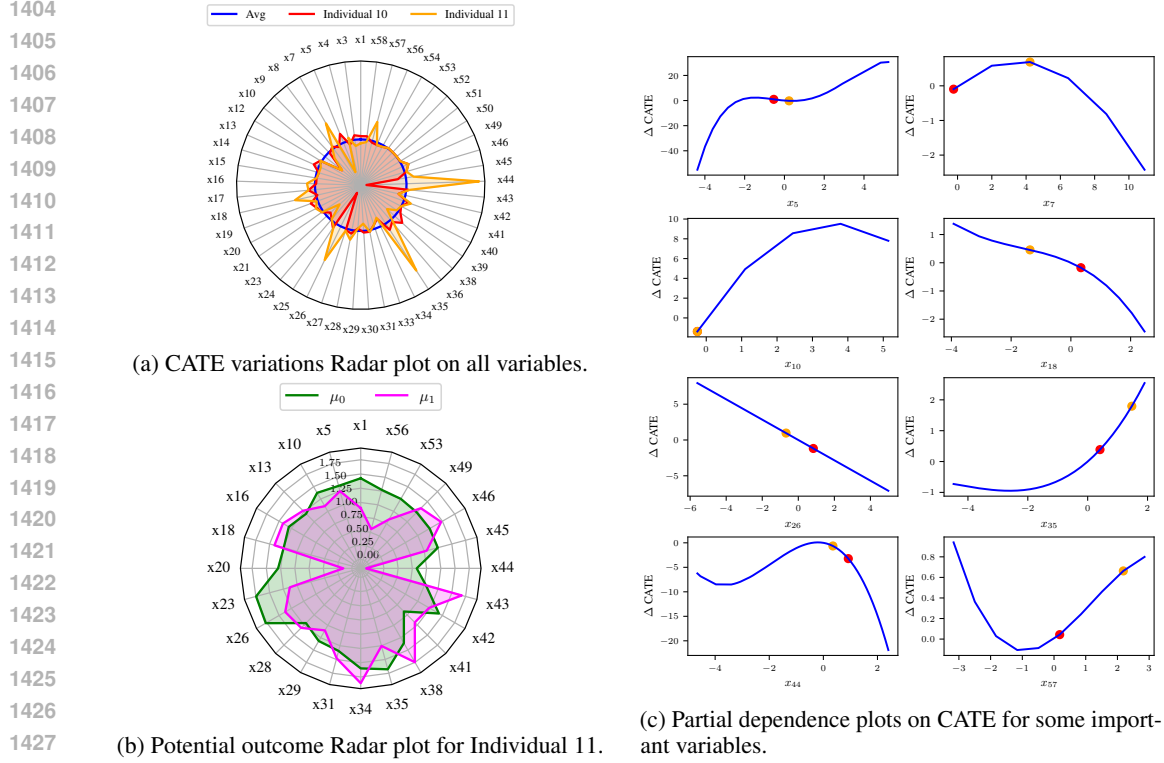
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

$$\begin{aligned}
\hat{\tau}(\mathbf{x}) = & 0.09 x_1^2 + 0.19 x_1 - 0.74 x_{10}^2 + 5.30 x_{10} + 0.01 x_{12}^3 - 0.05 x_{12}^2 - 0.09 x_{12} \\
& + 0.03 x_{13}^3 - 0.13 x_{13}^2 - 0.21 x_{13} - 0.06 x_{14}^2 + 0.33 x_{14} - 0.06 x_{15}^2 + 0.13 x_{15} \\
& + 0.19 x_{16}^2 - 0.22 x_{16} - 0.59 x_{17} - 0.03 x_{18}^3 - 0.14 x_{18}^2 - 0.47 x_{18} \\
& + 0.03 x_{19}^4 - 0.10 x_{19}^3 - 0.33 x_{19}^2 + 0.42 x_{19} + 0.14 x_{20}^3 - 0.29 x_{20}^2 - 0.30 x_{20} \\
& + 0.01 x_{21}^3 + 0.05 x_{21}^2 + 0.28 x_{21} + 0.25 x_{23}^3 - 0.50 x_{23}^2 - 0.83 x_{23} - 0.16 x_{24} \\
& - 0.07 x_{25}^2 + 0.14 x_{25} - 1.42 x_{26} - 0.01 x_{27}^4 + 0.04 x_{27}^3 - 0.06 x_{27}^2 + 0.06 x_{27} \\
& - 0.12 x_{28}^2 - 0.32 x_{28} + 0.00 x_{29}^4 + 0.08 x_{29}^3 + 0.07 x_{29}^2 + 0.01 x_{29} \\
& - 0.16 x_3 - 0.17 x_{30} - 0.01 x_{31}^2 - 0.08 x_{31} + 0.54 x_{33} \\
& - 0.03 x_{34}^2 - 0.17 x_{34} + 0.02 x_{35}^3 + 0.23 x_{35}^2 + 0.84 x_{35} \\
& + 0.03 x_{36}^3 - 0.14 x_{36}^2 + 0.22 x_{36} - 0.30 x_{38} - 0.15 x_{39}^2 - 0.20 x_{39} \\
& - 0.03 x_4^2 + 0.12 x_4 - 0.13 x_{40}^3 - 0.31 x_{40}^2 - 0.18 x_{40} \\
& + 0.04 x_{41}^3 - 0.39 x_{41}^2 + 0.02 x_{41} + 0.25 x_{42}^2 + 0.18 x_{42} \\
& + 0.01 x_{43}^4 - 0.03 x_{43}^3 - 0.09 x_{43}^2 + 0.08 x_{43} \\
& - 0.41 x_{44}^3 - 2.42 x_{44}^2 - 0.92 x_{44} - 0.06 x_{45}^3 - 0.29 x_{45}^2 - 0.32 x_{45} \\
& - 0.09 x_{46}^2 - 0.04 x_{46} - 1.10 x_{49} - 0.10 x_5^4 + 0.56 x_5^3 + 1.27 x_5^2 - 1.30 x_5 \\
& - 0.02 x_{50}^2 + 0.01 x_{50} - 0.12 x_{51} + 0.14 x_{52}^2 - 0.71 x_{52} - 0.57 x_{53} + 0.42 x_{54} \\
& - 0.02 x_{56}^3 + 0.01 x_{56}^2 + 0.52 x_{56} - 0.02 x_{57}^3 + 0.09 x_{57}^2 + 0.23 x_{57} \\
& - 0.02 x_{58}^3 - 0.17 x_{58}^2 - 0.34 x_{58} - 0.06 x_7^2 + 0.41 x_7 - 0.21 x_8 + 0.21 x_9^2 - 0.25 x_9 \\
& + 9.47.
\end{aligned}$$

This formula is long, due to the high number of covariates. Therefore, although the contribution of each variable has a polynomial form, analyzing the formula alone could be tricky. To help to gain intuitions about the contribution of each variable to the CATE variation, we present three useful plots in Fig. 13. In the plots, we can obtain the contributions of each feature to the causal effect. However, those visualizations should not be used to intervene in any feature except the treatment, since the variations presented in Fig. 13 cannot be interpreted as causal effect following our assumptions.

Lastly, we want to show that our *inductive bias* for simplicity of the atom selection is useful to improve interpretability. The formula below is provided by the standard auto-symbolic function (native of pyKAN (Liu et al., 2024b)). We observe polynomial relations with some covariates, but the nonlinear complex functions that follow are less interpretable than simple polynomials.

$$\begin{aligned}
\hat{\tau}(\mathbf{x}) = & -0.59 x_{17} - 0.16 x_{20} + 0.26 x_{21} - 1.41 x_{26} + 0.10 x_{27} + 0.28 x_{29} \\
& - 0.25 x_{31} + 0.54 x_{33} + 0.15 x_{34} + 0.34 x_{35} + 0.22 x_{36} - 0.30 x_{38} \\
& - 0.32 x_{40} + 0.02 x_{43} - 1.67 x_{44} - 0.40 x_{45} + 0.65 x_5 - 0.12 x_{51} \\
& - 0.57 x_{53} + 0.42 x_{54} - 0.02 x_{55} - 0.40 x_{58} \\
& - 0.00 (1.17 - 9.05 x_{50})^2 + 0.01 (3.95 - 3.98 x_{16})^2 \\
& - 0.00 (5.62 - 7.80 x_{41})^2 - 0.19 (6.55 - 1.78 x_{10})^2 \\
& - 0.02 (7.66 - 2.42 x_{10})^2 + 0.01 (9.47 - 3.82 x_{52})^2 \\
& + 0.00 (9.80 - 4.97 x_{46})^2 + 0.00 (-2.05 x_{16} - 5.95)^2 \\
& - 0.00 (-9.30 x_{23} - 2.38)^2 + 0.12 \sqrt{5.29 x_{31} + 2.54} \\
& - 0.00 (-8.38 x_{39} - 5.79)^2 - 1.03 \exp(0.60 x_{18}) - 0.23 \exp(0.97 x_{24}) \\
& - 0.53 \exp(0.56 x_{34}) + 0.47 \exp(0.87 x_{35}) + 1.12 \sin(0.63 x_1 - 7.20) \\
& - 0.40 \sin(0.81 x_1 - 0.81) - 0.24 \sin(4.88 x_{12} + 1.30) \\
& - 0.40 \sin(0.83 x_{13} - 1.02) + 1.09 \sin(5.19 x_{13} + 5.19) \\
& - 0.33 \sin(1.81 x_{14} + 2.36) + 2.42 \sin(9.01 x_{15} + 7.79) \\
& - 0.56 \sin(1.35 x_{19} - 2.19) - 0.26 \sin(1.14 x_{20} + 1.39)
\end{aligned}$$



1429 Figure 13: Visualization plots employing T-KAAM with ACIC-7. **(a)** Radar plot of the contribution of each variable to the CATE. In blue, the average predicted CATE in all the dataset. In red and orange, the respective contributions of each variable of individuals 10 and 11. **(b)** Radar plot of the potential outcomes of the individual 10. We can observe the contribution of each variable to the potential outcomes. Only the variables that have not been pruned in both subnetworks have been added to the plot. **(c)** Partial dependence plots of the CATE variation, given variations on the most important variables in the PRP. Particular values of CATE variation for the individuals 10 and 11 shows their contribution.

1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448

$$\begin{aligned}
& + 0.51 \sin(1.80 x_{23} - 10.00) - 0.50 \sin(0.63 x_{25} - 8.42) \\
& - 0.86 \sin(0.48 x_{28} - 4.36) + 2.25 \sin(0.52 x_{28} + 1.99) \\
& + 1.09 \sin(0.22 x_3 + 8.43) - 1.02 \sin(4.20 x_{30} + 1.41) \\
& + 0.58 \sin(7.77 x_{41} - 3.80) + 0.74 \sin(1.04 x_{42} - 8.42) \\
& + 0.75 \sin(5.35 x_{42} + 0.85) + 1.29 \sin(9.59 x_{46} + 3.61) \\
& + 0.53 \sin(4.58 x_{49} + 2.43) - 4.78 \sin(4.79 x_{49} - 0.60) \\
& + 0.88 \sin(4.21 x_{56} + 7.78) - 1.94 \sin(5.59 x_{56} - 7.58) \\
& + 0.41 \sin(0.88 x_{57} - 7.16) + 1.03 \sin(3.20 x_7 + 7.20) \\
& + 0.72 \sin(0.40 x_8 + 2.00) + 1.07 \sin(9.06 x_9 + 1.76) \\
& + 0.53 \tanh(0.98 x_{18} - 1.40) + 7.14 - 0.11 \exp(-0.93 x_4).
\end{aligned}$$

1449 C.3.2 HOMOGENEOUS CATE WITH S-KAAM

1450
1451 In this section, we illustrate how we obtain the *homogeneous* CATE (or, equivalently, the ATE) in the IHDP A dataset, where the causal effect of the treatment is known to be homogeneous and linear.

1452
1453 Therefore, we instantiate an S-KAAM, which yields the following formula in the potential outcome estimation.

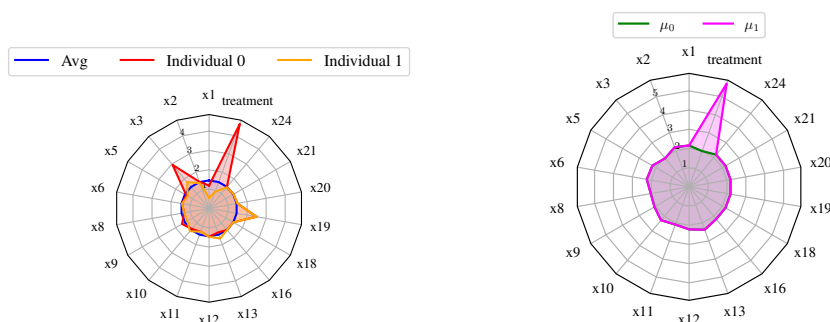
1454
1455
1456
1457

$$\begin{aligned}
\hat{\mu}(\mathbf{x}, \mathbf{t}) = & \underline{3.74} \mathbf{t} + 0.53 x_1 + 0.16 x_{10} + 0.59 x_{11} - 0.11 x_{12} + 0.34 x_{13} + 0.11 x_{16} - 0.17 x_{18} + 1.28 x_{19} \\
& + 0.01 x_2^3 - 0.01 x_2^2 - 0.23 x_{20} - 0.08 x_{21} + 0.11 x_{24} + 0.18 x_3^2 + 1.29 x_3 + 0.28 x_5
\end{aligned}$$

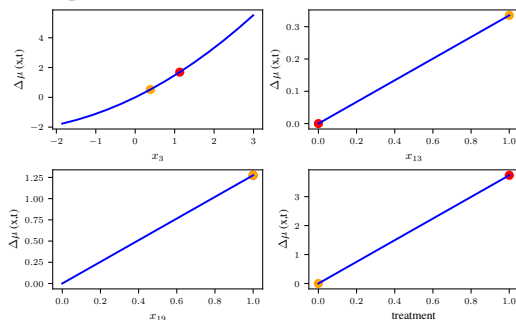
$$+ 0.01 x_6^4 - 0.03 x_6^3 + 0.03 x_6 + 1.47 x_8 + 0.32 x_9 + 1.88 .$$

As presented in Eq. 23, the CATE can be computed exclusively with the terms relative to t . In this case, the CATE can be directly extracted from the formula: 3.74.

We also represent in Fig. 14 some visualizations that we find interesting. First, for two given individuals, we present a PRP with the contribution of each variable to the variation of the predicted outcome, $\hat{\mu}(x, t)$, compared with the average of the predicted outcomes, $\mathbb{E}_{x, t}[\hat{\mu}(x, t)]$. On the right, we present the predicted potential outcomes for a given individual. As the effect is homogeneous (does not depend on the covariates), the contribution of each feature for both potential outcomes is the same, and the only difference is the causal effect of the treatment. Lastly, we present PDPs for treatment and other three variables (based on the radar plot), which present the variations of the predicted outcome with each variable. Following our assumptions, only the treatment curve can be seen as a causal effect, while the other curves can be used only to gain intuition of outcome behavior.



(a) Outcome variations Radar plot on all variables. (b) Potential outcome Radar plot for Individual 0.



(c) Partial dependence plots on CATE for some important variables.

Figure 14: Visualization plots employing S-KAAM with IHDP A. (a) Radar plot of the contribution of each variable to the outcome. In blue, the average predicted outcome in all the dataset. In red and orange, the respective contributions of each variable of individuals 0 and 1. (b) Radar plot of the potential outcomes of the individual 0. We can observe the contribution of each variable to the potential outcomes. (c) Partial dependence plots of the outcome variation, given variations on the most important variables in the PRP. Particularized for the individuals 0 and 1.

In this case, the standard symbolic substitution provides a similar formula as our proposal, due to the linearity of the dataset, and we omit its expression.

C.3.3 METRIC VARIATIONS FOLLOWING THE PIPELINE

We also show how the metrics—both observed (test loss) and unobserved (PEHE) in real data—vary when we perform the pruning and the symbolic substitution in the examples that we develop in App. C.3.

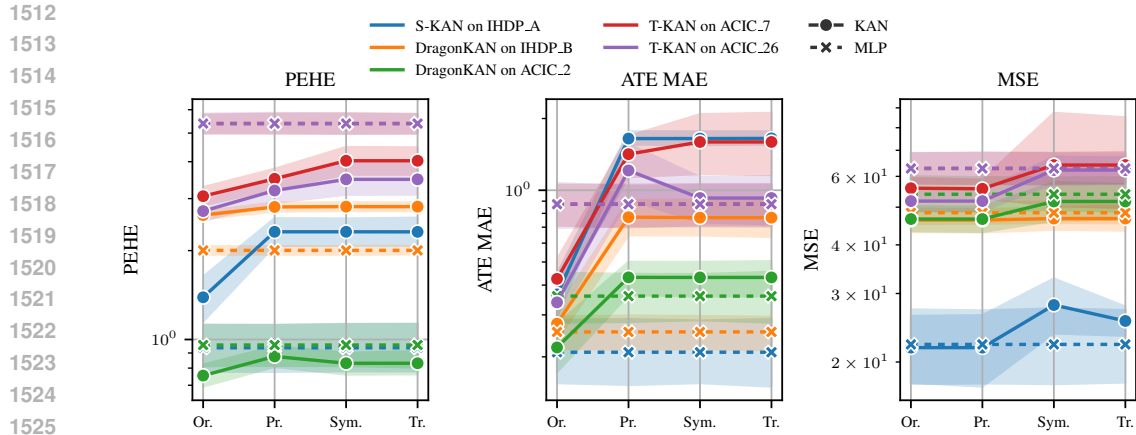


Figure 15: Variation of the metrics in each step of the pipeline: original (Or.), pruned (Pr.), Formula (For.) and 2-decimal truncation (Tr.).

We can observe in Fig. 15, for each dataset (with its respective model), the variation in performance when **i**) we prune the network, **ii**) we substitute the splines by symbolic formulas and **iii**) we truncate the formulas to have 2 decimals.

The conclusion of this experiment is that the metric variation in the different steps is high, so a practitioner should be careful when setting the budgets of step acceptance. MSE still shows signs of being a good proxy of the PEHE, representing the variations of the PEHE relatively well.

Note that, when following the pipeline, large errors can occur (e.g. when pruning S-KAN). A practitioner would either adjust the pruning threshold to limit changes in the estimated ITEs or simply reject pruning when it induces excessive error. For these experiments, we used the default pruning and symbolic substitution criteria ($\Gamma = 3 \cdot 10^{-2}$, $\Gamma_{R^2} = 0.98$) because the goal is to evaluate pruning uniformly across all benchmark datasets; tuning it per-dataset would create an unfair comparison.

C.4 EXPERIMENTS OF EXPRESSION RECOVERY

In this section, we add details about the experiments of CATE recovery, explained in §5.3. There, we have shown that for some expressions that can be captured by S-KAAM and T-KAAM respectively, the pipeline of causalKANs achieves better approximations of the true data-generating functions, than MLP or KANs without simplification steps. Those steps act as inductive biases towards a simpler representation.

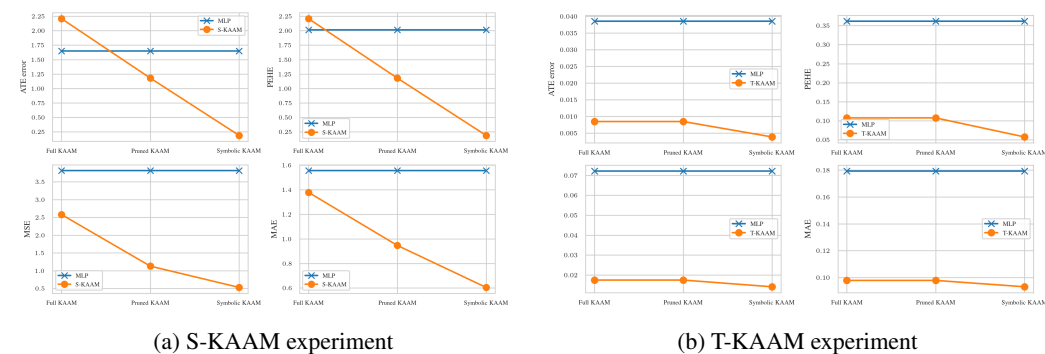


Figure 16: Metrics of the synthetic experiments that show function recovery: ATE error, PEHE, MSE and MAE. ATE error and PEHE can only be computed observing the counterfactual, while MAE and MSE are functions of observational data. **Lower is better.**

We can observe in Fig. 16 that all metrics decrease (improve) while computing the simplification steps, specially the symbolic substitution. That information is complementary to the curves observed in Fig. 6, that show a better fit of the symbolic curves. After symbolic substitution, the PEHE and the ATE error are much lower than the achieved by the causalNN counterpart in each case.

C.5 TIME CONSUMPTION

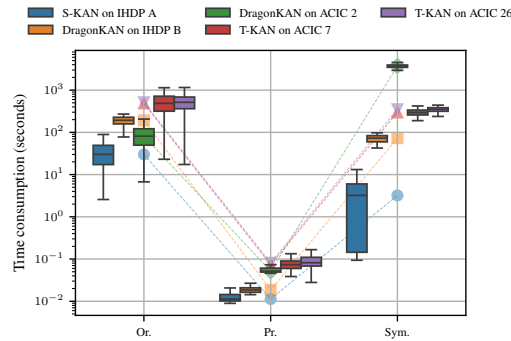
We are interested in comparing training and inference time of causalKANs with their respective causalNNs. It is well known that KANs take more time to train than MLPs, for networks with the same number of parameters. However, KANs usually require less parameters than MLPs to achieve similar performance (Liu et al., 2024b). For the tasks analyzed in this paper, causalKANs requires greater computational effort during both training and inference than existing causal neural networks. We consider this additional cost justified by the interpretability that causalKANs provides, which we regard as a central advantage for causal analysis. We report the relative training and inference times in Tab 2.

Dataset	Model Name	KAN		MLP		Ratio (Train)
		Training (s)	Inference (s)	Training (s)	Inference (s)	KAN/MLP
IHDP A	S-Learner	23.62 _{19.25}	0.12 _{0.02}	1.98 _{2.05}	0.00 _{0.00}	12
	T-Learner	97.89 _{75.57}	0.03 _{0.01}	29.00 _{29.19}	0.00 _{0.00}	7/2
	TarNet	67.21 _{52.39}	0.01 _{0.00}	8.02 _{9.27}	0.00 _{0.00}	8
	DragonNet	62.92 _{48.39}	0.01 _{0.00}	9.74 _{10.28}	0.00 _{0.00}	6/1
IHDP B	S-Learner	56.24 _{22.13}	0.06 _{0.01}	2.61 _{1.24}	0.00 _{0.00}	22
	T-Learner	58.49 _{14.10}	0.03 _{0.00}	45.81 _{18.08}	0.00 _{0.00}	9/7
	TarNet	31.39 _{10.88}	0.03 _{0.00}	42.37 _{16.43}	0.00 _{0.00}	3/4
	DragonNet	157.82 _{38.35}	0.01 _{0.00}	187.26 _{31.89}	0.00 _{0.00}	5/6
ACIC 2	S-Learner	120.96 _{54.54}	0.05 _{0.01}	23.13 _{15.77}	0.00 _{0.00}	5
	T-Learner	368.94 _{153.74}	0.22 _{0.03}	162.59 _{28.60}	0.01 _{0.00}	9/4
	TarNet	135.38 _{83.87}	0.05 _{0.01}	55.08 _{62.84}	0.01 _{0.00}	5/2
	DragonNet	92.83 _{55.77}	0.04 _{0.01}	63.17 _{70.06}	0.01 _{0.00}	3/2
ACIC 7	S-Learner	147.31 _{110.57}	0.23 _{0.03}	4.56 _{7.80}	0.00 _{0.00}	32
	T-Learner	377.66 _{177.18}	0.04 _{0.01}	39.90 _{43.24}	0.01 _{0.00}	9
	TarNet	236.88 _{109.75}	0.03 _{0.00}	40.38 _{50.22}	0.01 _{0.00}	6
	DragonNet	272.19 _{127.14}	0.03 _{0.00}	51.86 _{65.79}	0.01 _{0.00}	5
ACIC 26	S-Learner	141.11 _{96.81}	0.17 _{0.02}	4.46 _{7.54}	0.00 _{0.00}	32
	T-Learner	145.50 _{70.85}	0.03 _{0.00}	39.36 _{42.07}	0.01 _{0.00}	7/2
	TarNet	111.13 _{53.19}	0.02 _{0.00}	39.89 _{49.73}	0.01 _{0.00}	8/3
	DragonNet	147.49 _{70.45}	0.03 _{0.05}	48.26 _{58.27}	0.01 _{0.00}	3
NSLM	S-Learner	75.62 _{44.37}	0.38 _{0.09}	73.57 _{29.65}	0.01 _{0.00}	1
	T-Learner	12.14 _{3.92}	0.06 _{0.01}	155.73 _{89.40}	0.01 _{0.00}	1/13
	TarNet	10.70 _{2.46}	0.07 _{0.01}	130.49 _{63.07}	0.00 _{0.00}	1/12
	DragonNet	19.61 _{3.92}	0.06 _{0.01}	214.24 _{34.63}	0.00 _{0.00}	1/11
NEWS	S-Learner	40.28 _{18.32}	0.24 _{0.04}	2.53 _{0.43}	0.01 _{0.00}	16
	T-Learner	227.39 _{94.20}	0.23 _{0.28}	147.36 _{77.26}	0.00 _{0.00}	3/2
	TarNet	45.75 _{17.19}	0.08 _{0.17}	5.87 _{0.50}	0.00 _{0.00}	8
	DragonNet	19.32 _{5.97}	0.06 _{0.01}	6.54 _{0.41}	0.00 _{0.00}	3
TCGA	S-Learner	122.09 _{37.81}	0.31 _{0.23}	48.05 _{18.54}	0.01 _{0.00}	5/2
	T-Learner	269.33 _{101.08}	0.19 _{0.12}	778.61 _{338.08}	0.00 _{0.01}	1/3
	TarNet	285.42 _{117.24}	0.10 _{0.06}	246.19 _{109.33}	0.01 _{0.00}	6/5
	DragonNet	185.98 _{100.10}	0.09 _{0.02}	31.35 _{2.40}	0.01 _{0.00}	6

Table 2: Comparison of training and inference times (in seconds) for KAN and MLP across datasets. The last column shows the ratio of KAN to MLP training time, approximated as simple fractions. Values are reported as mean_{std}.

In average, causalKANs train slower than causalNNs: KAN training time $\approx 8 \times$ MLP training time. However, note that the time consumption depends largely on the dataset, since we can observe that, for the NSLM dataset, causalKANs trains much faster than MLP.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632



1633 Figure 17: Training of original network (Or.), pruning (Pr.) and symbolic (Sym.) time consumption in seconds for IHDP and ACIC datasets. Logarithmic scale. Boxplots represent the distribution of times in all realizations. Markers and lines represent the medians of each distribution.

1636
1637
1638
1639
1640
1641
1642
1643

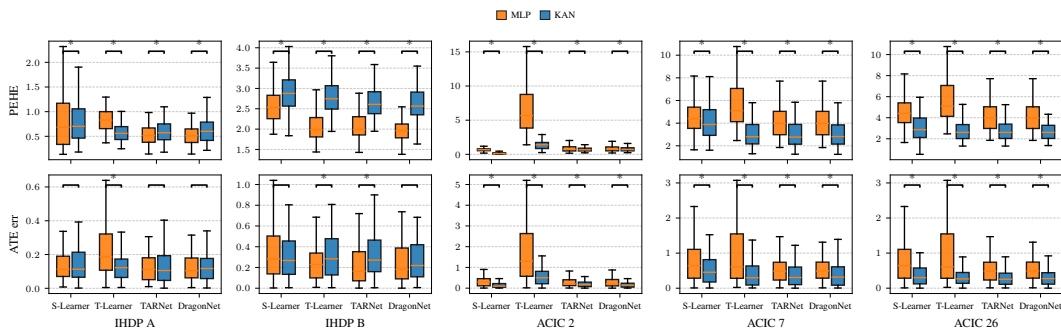
Pipeline time consumption. Besides training, the rest of the interpretability pipeline (Pruning and symbolic substitution) also introduces complexity in terms of computation. We include in Fig. 17 the distribution of training (Or.), pruning (Pr.) and symbolic substitution (Sym.) times for the best causalKANs in the ACIC and IHDP datasets. All these training times have been measured in a CPU AMD Ryzen Threadripper 7970X 32-Cores. We can observe that the symbolic substitution time is comparable or even higher than training time in some datasets. On the other hand, pruning time is negligible.

1644 C.6 COMPARISON VISUALIZATIONS

1645

In the same fashion, we have compared each causalNN with its respective causalKAN, in Fig. 18.

1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658



1659
1660
1661
1662
1663

Figure 18: Comparison of KAN vs MLP across datasets. Top row: PEHE, bottom row: ATE err. The trend of the difference is not constant across datasets. For example, we can observe that causalKANs achieve better PEHE metrics in IHDP A and ACIC 2/7/26, but worse in IHDP B, than their respective causalNNs.

1664
1665
1666

We also include, in Tab 3, p-values for transparency, as they provide a more nuanced understanding of the evidence against the null hypothesis than a binary significant/non-significant determination at a fixed α level.

1667 D COMPLETE PIPELINE

1668

1669
1670
1671
1672
1673

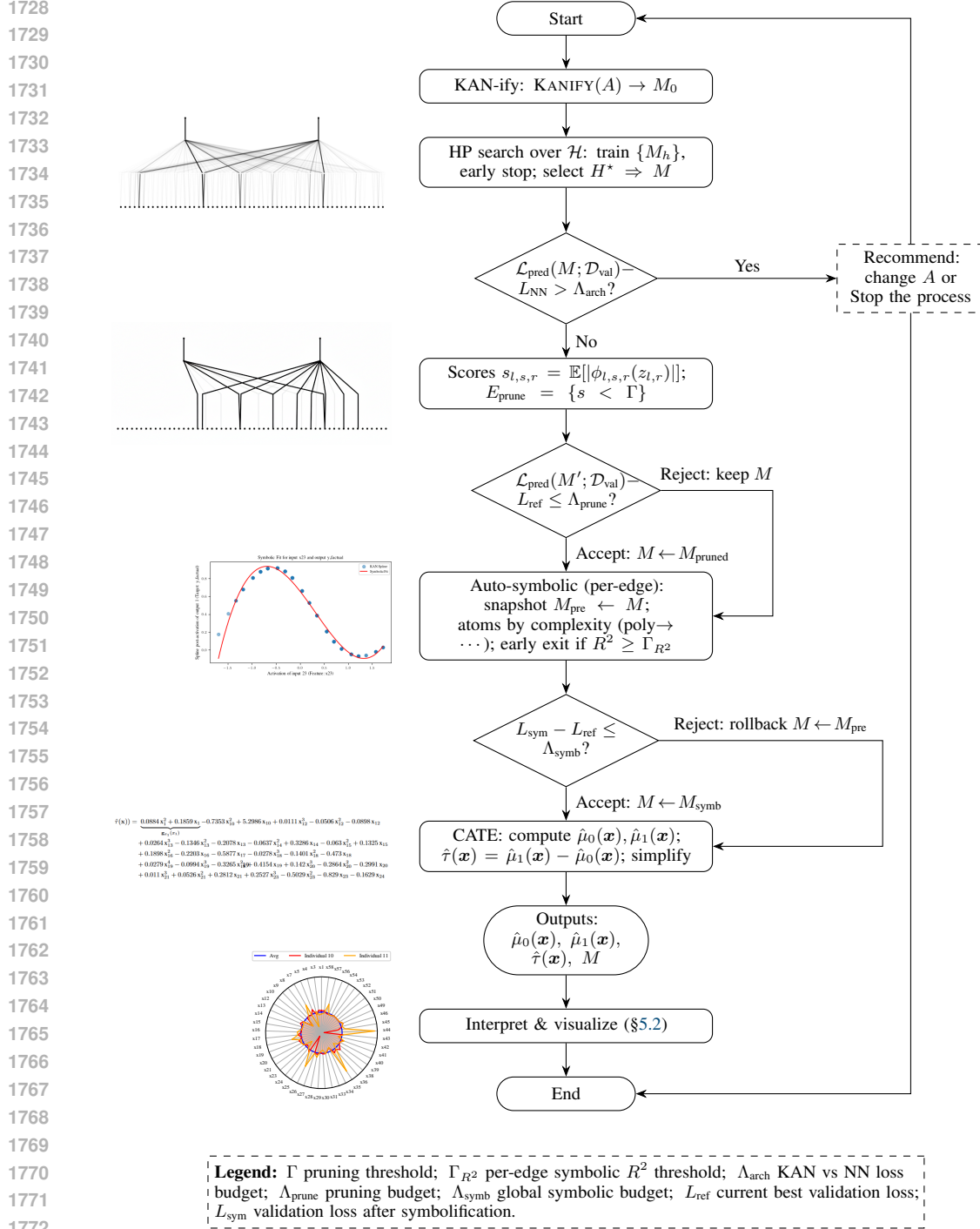
We offer two different complete visualizations of the pipeline that we propose, in an algorithm Alg. 1 version that details all the steps, including the details and the computational step of each block, and a block diagram Fig. 19 in which is easier to focus on the decision steps and the importance of the budgets.

Dataset	Architecture	KAN		MLP	
		p (ATE err)	p (PEHE)	p (ATE err)	p (PEHE)
IHDP A	S-Learner	0.781	0.121	0.132	$< 10^{-3}$
	T-Learner	0.889	1.000	$< 10^{-3}$	$< 10^{-3}$
	TarNet	1.000	1.000	0.906	1.000
	DragonNet	0.906	1.000	1.000	1.000
IHDP B	S-Learner	0.006	$< 10^{-3}$	0.012	$< 10^{-3}$
	T-Learner	0.007	$< 10^{-3}$	0.388	0.189
	TarNet	0.003	$< 10^{-3}$	1.000	0.189
	DragonNet	0.388	0.006	0.436	1.000
ACIC 2	S-Learner	1.000	1.000	$< 10^{-3}$	$< 10^{-3}$
	T-Learner	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$
	TarNet	0.013	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$
	DragonNet	0.157	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$
ACIC 7	S-Learner	0.020	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$
	T-Learner	0.908	0.788	$< 10^{-3}$	$< 10^{-3}$
	TarNet	1.000	0.788	0.020	$< 10^{-3}$
	DragonNet	0.706	1.000	$< 10^{-3}$	$< 10^{-3}$
ACIC 26	S-Learner	0.838	0.022	$< 10^{-3}$	$< 10^{-3}$
	T-Learner	1.000	0.511	$< 10^{-3}$	$< 10^{-3}$
	TarNet	1.000	0.511	0.011	$< 10^{-3}$
	DragonNet	1.000	1.000	$< 10^{-3}$	$< 10^{-3}$
NSLM	S-Learner	1.000	1.000	$< 10^{-3}$	$< 10^{-3}$
	T-Learner	0.15	1.000	0.52	$< 10^{-3}$
	TarNet	0.01	1.000	0.44	$< 10^{-3}$
	DragonNet	$< 10^{-3}$	0.104	0.196	$< 10^{-3}$
NEWS	S-Learner	0.621	0.406	$< 10^{-3}$	$< 10^{-3}$
	T-Learner	0.506	$< 10^{-3}$	$< 10^{-3}$	0.006
	TarNet	1.000	0.006	0.101	1.000
	DragonNet	0.621	0.003	$< 10^{-3}$	0.488
TCGA	S-Learner	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$	0.039
	T-Learner	$< 10^{-3}$	$< 10^{-3}$	1.000	$< 10^{-3}$
	TarNet	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$	1.000
	DragonNet	$< 10^{-3}$	$< 10^{-3}$	$< 10^{-3}$	0.001

Table 3: Friedman post-hoc *corrected* p -values per dataset and metric, arranged to mirror Tab 1. Baselines are underlined (reported as $p = 1.000$) and methods not significantly different from the baseline at $\alpha = 0.05$ are in **bold**. Very small p -values are reported as $< 10^{-3}$.

As can be observed, although the pipeline is well defined, there are many variables (or hyperparameters) that the practitioner should vary depending on the dataset or the application of our proposal. We include here some recommendations to choose these parameters and some discussion about the pipeline.

First of all, note that we treat KANs as a modeling option, not a universally superior alternative. Therefore, the practitioner should decide in each step if the causalKAN should be used or if one should keep the causalNN baseline. For example, regarding the benchmarking experiments of §5, we observe that TCGA obtains ITEs that are very far from the causalNNs estimates. Therefore, the practitioner should evaluate if that is acceptable. The same criteria should be applied with the rest of the thresholds. For example, when pruning the DragonKAN in IHDP B in Fig. 15, we observe that both PEHE and ATE move away from causalNN estimates. A practitioner should establish a threshold of admissible deviation. However, note that the variation in metrics given by the pipeline steps are dependant of the parameters established by the methods applied: in the pruning step, the pruning threshold Γ ; in the symbolic substitution, the Γ_{R^2} . We recommend to test several of these parameters, and compare the deviation with the pipeline budgets (Λ). Therefore, the process of selecting those parameters can be cyclic and the practitioner should establish the number of iterations that are admissible. For the experiments, we set the default parameters of pruning (edge threshold



1773 Figure 19: Block diagram of CausalKAN pipeline with explicit Accept/Reject semantics and final
1774 outputs.

1775
1776 $3 \cdot 10^{-2}$, node threshold 10^{-2}) and R^2 ($\Gamma_{R^2}=0.98$) for getting a fair comparison. Therefore, the
1777 accept-reject gates are not static steps, but should be defined depending on the dataset and task.
1778
1779
1780
1781

1782 E LLM USAGE
1783

1784 Large language models (LLMs) were used exclusively to improve the clarity and grammar of the
1785 manuscript, and to assist with the presentation of plots and figures in Python and \LaTeX (TikZ). All
1786 scientific contributions, methodological developments, experiments, and analyses were conceived,
1787 implemented, and validated by the authors, with no use of LLMs.
1788

1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

1836 **Algorithm 1** CausalKAN : Interpretable CATE via KAN-ified causal networks

1837 **Input:** Base causalNN A ; splits $\mathcal{D}_{\text{tr}}, \mathcal{D}_{\text{val}}, \mathcal{D}_{\text{te}}$ with $(\mathbf{x}, \mathbf{t}, \mathbf{y})$; budgets $\Lambda_{\text{prune}}, \Lambda_{\text{symp}}$; thresholds Γ (prune), Γ_{R^2}
1838 (symbolic), Λ_{arch} (KAN vs NN); HP spaces, \mathcal{H} , (depth, width, spline grid, $\lambda_1, \lambda_c, \lambda_s, \lambda_H$); atom dict
1839 $\{f_m\}_{m=1}^M$ ordered by complexity (poly \rightarrow trigs \rightarrow others); optimizer, early stopping.

1840 **Output:** $\hat{\mu}_0(\mathbf{x}), \hat{\mu}_1(\mathbf{x}), \hat{\tau}(\mathbf{x}) \equiv \hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$

1841

1842 1: **KAN-ification:** $M_0 \leftarrow \text{KANIFY}(A)$

1843

1844 2: **HP search & training:**

1845 3: **for** $h \in \mathcal{H}$ **do**

1846 4: $M_h \leftarrow \text{INSTANTIATE}(M_0, h)$

1847 5: Minimize on \mathcal{D}_{tr} :

1848
$$\mathcal{L} = \mathcal{L}_{\text{pred}}(M_h) + \lambda_1 \sum_{l,s,r} \mathbb{E}[|\phi_{l,s,r}(z_{l,r})|] + \lambda_c \sum |w_s| + \lambda_s \sum |w_s - w'_s| + \lambda_H \sum H$$

1849

1850 6: Early stop on \mathcal{D}_{val} ; store $L_{\text{val}}(h) = \mathcal{L}_{\text{pred}}(M_h^\dagger; \mathcal{D}_{\text{val}})$

1851 7: **end for**

1852 8: $H^* \leftarrow \arg \min_h L_{\text{val}}(h)$ with tie-break by simplicity (fewer layers, smaller grids/nodes, no MultKAN)

1853 9: $M \leftarrow M_{H^*}^\dagger$; $L_{\text{ref}} \leftarrow L_{\text{val}}(H^*)$

1854

1855 10: **Baseline causalNN check (KAN vs NN):**

1856 11: Train A^\dagger (original causalNN) under its best HPs on \mathcal{D}_{tr} with early stopping

1857 12: $L_{\text{NN}} \leftarrow \mathcal{L}_{\text{pred}}(A^\dagger; \mathcal{D}_{\text{val}})$

1858 13: **if** $\mathcal{L}_{\text{pred}}(M; \mathcal{D}_{\text{val}}) - L_{\text{NN}} > \Lambda_{\text{arch}}$ **then**

1859 14: **Warn/Recommend:** change A (architecture underperforming);

1860 15: **end if**

1861

1862 16: **Pruning (accept-reject):**

1863 17: $s_{l,s,r} \leftarrow \mathbb{E}_{\mathcal{D}_{\text{val}}} [|\phi_{l,s,r}(z_{l,r})|]$

1864 18: $E_{\text{prune}} \leftarrow \{(l, s \rightarrow r) : s_{l,s,r} < \Gamma\}$

1865 19: **if** $E_{\text{prune}} \neq \emptyset$ **then**

1866 20: $M' \leftarrow \text{REMOVEEDGESANDISOLATEDNODES}(M, E_{\text{prune}})$

1867 21: **if** $\mathcal{L}_{\text{pred}}(M'; \mathcal{D}_{\text{val}}) - L_{\text{ref}} \leq \Lambda_{\text{prune}}$ **then**

1868 22: $M \leftarrow M'$; $L_{\text{ref}} \leftarrow \mathcal{L}_{\text{pred}}(M; \mathcal{D}_{\text{val}})$

1869 23: **else**

1870 24: **Reject pruning**

1871 25: **end if**

1872 26: **end if**

1873

1874 27: **Auto-symbolic (per-edge, early exit by R^2) + global accept:**

1875 28: $M_{\text{pre}} \leftarrow M$ ▷ snapshot for possible rollback

1876 29: **for** each edge $e = (l, s \rightarrow r)$ in M **do**

1877 30: collect $\{(u_k, v_k)\}$ on \mathcal{D}_{val} : $u_k = z_{l,r}^{(k)}$, $v_k = \phi_{l,s,r}(u_k)$

1878 31: **for** $m = 1 \rightarrow M$ (**complexity-ordered**) **do**

1879 32: $(a^*, b^*, c^*, d^*) \leftarrow \arg \min_{a,b,c,d} \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_k (v_k - [cf_m(au_k + b) + d])^2$

1880 33: compute R_m^2 on \mathcal{D}_{val}

1881 34: **if** $R_m^2 \geq \Gamma_{R^2}$ **then**

1882 35: $M \leftarrow \text{REPLACEEDGEWITHATOM}(M, e, f_m, a^*, b^*, c^*, d^*)$ ▷ early accept for this edge

1883 36: **break**

1884 37: **end if**

1885 38: **end for**

1886 39: **end for**

1887 40: $L_{\text{sym}} \leftarrow \mathcal{L}_{\text{pred}}(M; \mathcal{D}_{\text{val}})$

1888 41: **if** $L_{\text{sym}} - L_{\text{ref}} \leq \Lambda_{\text{symp}}$ **then**

1889 42: **Accept** symbolic model; $L_{\text{ref}} \leftarrow L_{\text{sym}}$

1890 43: **else**

1891 44: **Reject** symbolic model; $M \leftarrow M_{\text{pre}}$

1892 45: **end if**

1893

1894 46: **CATE extraction (difference \Rightarrow simplify):**

1895 47: Compute $\hat{\mu}_0(\mathbf{x}), \hat{\mu}_1(\mathbf{x})$ by forward evaluation of the two heads

1896 48: $\hat{\tau}(\mathbf{x}) \leftarrow \hat{\mu}_1(\mathbf{x}) - \hat{\mu}_0(\mathbf{x})$

1897 49: $\text{SIMPLIFYALGEBRA}(\hat{\tau}(\mathbf{x}))$ ▷ cancel/factor common terms *after* the difference

1898 50: **return** $(\hat{\mu}_0, \hat{\mu}_1, \hat{\tau}(\mathbf{x}), M)$
