ANOVA-NODE: AN IDENTIFIABLE NEURAL NET WORK FOR THE FUNCTIONAL ANOVA MODEL FOR BETTER INTERPRETABILITY

Anonymous authors Paper under double-blind review

ABSTRACT

Interpretability for machine learning models is becoming more and more important as machine learning models become more complex. The functional ANOVA model, which decomposes a high-dimensional function into a sum of lower dimensional functions so called *components*, is one of the most popular tools for interpretable AI, and recently, various neural network models have been developed for estimating each component in the functional ANOVA model. However, such neural networks are highly unstable when estimating components since the components themselves are not uniquely defined. That is, there are multiple functional ANOVA decompositions for a given function. In this paper, we propose a novel interpretable model which guarantees a unique functional ANOVA decomposition and thus is able to estimate each component stably. We call our proposed model ANOVA-NODE since it is a modification of Neural Oblivious Decision Ensembles (NODE) for the functional ANOVA model. Theoretically, we prove that ANOVA-NODE can approximate a smooth function well. Additionally, we experimentally show that ANOVA-NODE provides much more stable estimation of each component and thus much more stable interpretation when training data and initial values of the model parameters vary than existing neural network models do.

006

007

012 013

014

015

016

017

018

019

021

022

024

025

026

027

028

1 INTRODUCTION

Interpretability has become more important as artificial intelligence (AI) models have become more sophisticated and complicated in recent years. Various methods for interpretable AI can be catego-rized into two groups. One is transparent box design, where interpretable machine learning models such as linear models and decision trees are used to learn a prediction model. Typically, interpretable
models are inferior to black box models in terms of prediction powers. The other group consists of post-hoc interpretation methods that try to interpret a given black box models (Lundberg, 2017; Ribeiro et al., 2016). While post-hoc interpretation methods do not hamper the prediction power of a given black box model at all, they often exhibit instability and lack faithfulness (Slack et al., 2020).

In this paper, we focus on the transparent box design based on the functional ANOVA model (Ho-effding, 1992). The functional ANOVA model approximates a given complex high-dimensional function by the sum of low dimensional (e.g., one or two dimensional) functions referred to as *components*. One of the most representative examples of the functional ANOVA model is the generalized additive model (GAM, Hastie & Tibshirani (1987)), which consists of the summation of one-dimensional functions, each corresponding to an input feature. Low dimensional functions are easier to understand, and thus the functional ANOVA model is popularly used for interpretable AI (Lengerich et al., 2020; Märtens & Yau, 2020).

Recently, various learning algorithms for the functional ANOVA model based on neural networks have been proposed (Agarwal et al., 2021; Radenovic et al., 2022; Chang et al., 2021). While they provide accurate prediction models, existing neural networks struggle to estimate each component in the functional ANOVA model due to unidentifiability (Lengerich et al., 2020). That is, completely different components could result in the same prediction model. Figure 1 presents the plots of the main effects estimated by NA²M, NB²M, and our proposed ANOVA-N²ODE on CALIFORNIA



Learning the functional ANOVA model by using neural networks under such an identifiability constraint, however, is not easy since a standard gradient descent algorithm could not be applicable.
The aim of this paper is to develop a specially designed neural network that automatically satisfies the sum-to-zero condition and thus can be learned by a standard gradient descent algorithm. The proposed neural network is a modification of Neural Oblivious Decision Ensembles (NODE, Popov et al. (2019)) for the functional ANOVA model. Chang et al. (2021) propose a modification of NODE for the GAM called NODE-GAM, but NODE-GAM is not effective in estimating each component due to the identifiability issue. We modify NODE such that each component in the functional ANOVA model is estimated uniquely. We call our algorithm ANOVA-NODE.

There exist various learning algorithms for the functional ANOVA model under the sum-to-zero condition such as Gu & Wahba (1993); Lin & Zhang (2006); Kim et al. (2009) that do not use neural networks. An advantage of ANOVA-NODE compared to these non-neural algorithms is that a gradient descent based optimization algorithm can be used in learning and hence end-to-end learning is possible when it is combined with other neural networks. See Section 4.5 for an example.

7 Overall, our contributions are as follows.

- We propose a novel XAI model (ANOVA-NODE) which trains the functional ANOVA model under the sum-to-zero condition using a gradient descent-based optimization algorithm.
- We prove the universal approximation property in the sense that ANOVA-NODE can approximate any smooth function up to an arbitrary precision.
- By analyzing multiple benchmark datasets, we illustrate that ANOVA-NODE provides more stable estimation and interpretation of each component compared to the baseline models, including NAM (Agarwal et al., 2021), NBM (Radenovic et al., 2022), NODE-GAM (Chang et al., 2021) and XGB (Chen & Guestrin, 2016) without losing prediction accuracy.

¹⁰⁸ 2 BACKGROUND

110 2.1 NOTATION

111 Let $\mathbf{x} = (x_1, ..., x_p)^\top \in \mathcal{X} = \mathcal{X}_1 \times ... \times \mathcal{X}_p$ be a vector of input features, where we assume 112 $\mathcal{X} \subseteq [-a, a]^p$ for some a > 0. We denote $[p] = \{1, ..., p\}$, and denote its power set as $\mathcal{P}([p])$. For 113 $\mathbf{x} \in \mathcal{X}$ and $S \subseteq [p]$, let $\mathbf{x}_S = (x_j, j \in S)^\top$. We denote f_S as a function of \mathbf{x}_S . For a real-valued 114 function $f : \mathcal{X} \to \mathbb{R}$, we denote $||f||_{\infty} = \sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x})|$.

115 2.2 FUNCTIONAL ANOVA MODEL

The functional ANOVA model (Hoeffding, 1992) decomposes a high-dimensional function f into the sum of low-dimensional functions

119 120

121

129 130

131

135 136

139

140

 $f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p f_j(x_j) + \sum_{j < k} f_{jk}(x_j, x_k) + \cdots,$

which is considered as one of the most important XAI tools. Here, $f(\mathbf{x}) = g(\mathbb{E}(Y|X = x))$ where g is a link function and Y is target variable, $f_j, j \in [p]$ are called the main effects, and $f_{j,k}, (j,k) \in [p]^2$ are called the second interaction terms and so on. In practice, only interactions of lower orders (e.g., the main and second order only) are included in the decomposition for a more transparent interpretation.

Generalized Additive Model (GAM, Hastie & Tibshirani (1987)) is a special form of the functional
 ANOVA model where only the main effects are included in the model, that is

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p f_j(x_j).$$

Similarly, GA^2M is defined as the functional ANOVA model including all of the main effects and second order interactions,

$$f(\mathbf{x}) = \beta_0 + \sum_{S \subseteq [p], |S| \le 2} f_S(\mathbf{x}_S)$$

or more generally we can consider GA^dM defined as

$$f(\mathbf{x}) = \beta_0 + \sum_{S \subseteq [p], |S| \le d} f_S(\mathbf{x}_S).$$

Several learning algorithms for the functional ANOVA model have been proposed. Gu & Wahba (1993) applied the smoothing spline to learn the functional ANOVA model, Lin & Zhang (2006) developed a component-wise sparse penalty, and Kim et al. (2009) proposed a boosting algorithm for the functional ANOVA model. In addition, the functional ANOVA model has been applied to various problems such as sensitivity analysis (Chastaing et al., 2012), survival analysis (Huang et al., 2000), diagnostics of high-dimensional functions (Hooker, 2007) and machine learning models (Lengerich et al., 2020; Märtens & Yau, 2020).

Recently, learning the functional ANOVA model using neural networks has received much attention since gradient descent based learning algorithms can be easily scaled up. Examples are Neural Additive Model (NAM, Agarwal et al. (2021)), Neural Basis Model (NBM, Radenovic et al. (2022)) and NODE-GAM (Chang et al., 2021). NAM uses deep neural network (DNN) to train each component of GAM. NBM achieves a significant reduction in training time compared to NAM by using basis DNNs to train all components. NODE-GAM extends Neural Obilvious Decision Ensembles (NODE, Popov et al. (2019)) for GAM.

155 156

2.3 NEURAL OBLIVIOUS DECISION ENSEMBLES (NODE)

Oblivious Decision Tree (ODT) (Kohavi & Li, 1995) is a decision tree which has the following two properties:

- 160
- 1. All terminal nodes have the same depth.
 - 2. All rules at each depth are identical.

162 Let d be the depth of ODT, for $t \leq d$, $F^t(\mathbf{x})$ be the feature function used for the rule at the depth 163 t, bt is the split value in the depth t and $\mathcal{B} \in \mathbb{R}^{2^d}$ be the height parameters at the terminal nodes of 164 ODT. Then, the ODT model is given as 165

$$h(\mathbf{x}) = \mathcal{B}'\left(\begin{bmatrix} \mathbb{I}(F^1(\mathbf{x}) \le b^1) \\ \mathbb{I}(F^1(\mathbf{x}) > b^1) \end{bmatrix} \otimes \begin{bmatrix} \mathbb{I}(F^2(\mathbf{x}) \le b^2) \\ \mathbb{I}(F^2(\mathbf{x}) > b^2) \end{bmatrix} \otimes \dots \otimes \begin{bmatrix} \mathbb{I}(F^d(\mathbf{x}) \le b^d) \\ \mathbb{I}(F^d(\mathbf{x}) > b^d) \end{bmatrix} \right)$$

where \mathbb{I} is the indicator function and \otimes is the outer product. 169

To train ODT using gradient descent methods, Popov et al. (2019) replaced the indicator function 170 $\mathbb{I}(F(\mathbf{x}) > b)$ with the $entmax_{\nu}((\frac{F(\mathbf{x})-b}{\gamma}, 0)')_1$ (Peters et al., 2019) which is differentiable. More-171 over, the feature function F^t is a weighed sum of input features. Note that $entmax_{\nu}((\frac{F(\mathbf{x})-b}{\gamma},0)')_1$ 172 173 works similarly to a sparse sigmoid. They referred to this ODT as Neural ODT (NODT) 174

In addition, Popov et al. (2019) proposed a layer architecture which involves the output of the current layer being concatenated with the current input and then fed into the next layer. The final output of NODT with a layer architecture can be obtained by averaging the outputs of each layer. Finally, they proposed Neural Oblivious Decision Ensembles (NODE) as an ensemble of NODTs.

178 179 180

183

185 186 187

188 189

199

201

213 214

175

176

177

166 167

> 3 PROPOSED MODEL

181 3.1 IDENTIFIABILITY ISSUE 182

An unsolved but important problem in neural functional ANOVA algorithms is the *identifiability* of each component. The functional ANOVA model itself is not identifiable. That is, there are multiple functional ANOVA decompositions of a given function. For example,

$$f(x_1, x_2) = f_1(x_1) + f_2(x_2) + f_{12}(x_1, x_2)$$

where $f_1(x_1) = x_1, f_2(x_2) = x_2, f_{1,2}(x_1, x_2) = x_1x_2$ can be expressed as

$$f(x_1, x_2) = f_1^*(x_1) + f_2^*(x_2) + f_{1,2}^*(x_1, x_2)$$

190 where $f_1^*(x_1) = -x_1, f_2^*(x_2) = x_2, f_{12}^*(x_1, x_2) = x_1(x_2 + 2)$ (Lengerich et al., 2020). Without the 191 identifiability, each component cannot be estimated uniquely and thus interpretation of the model 192 becomes unstable and inaccurate. 193

A simple remedy to ensure the identifiability of each component is to put constraints. One of the 194 most popular constraints for the identifiability of the components in the functional ANOVA model is so called the *sum-to-zero* condition (Gu & Wahba, 1993; Kim et al., 2009). When we consider the 196 functional ANOVA model using interactions in $\mathbf{S} \subseteq \mathcal{P}([p])$, the *sum-to-zero* condition is 197

$$\forall S \in \mathbf{S}, \ \forall j \in S, \ \forall \mathbf{z} \in \mathcal{X}_{S \setminus \{j\}}, \ \int_{\mathcal{X}_j} f_S(\mathbf{x}_{S \setminus \{j\}} = \mathbf{z}, x_j) \mu_j(dx_j) = 0 \tag{1}$$

200 for some probability measure μ_i on \mathcal{X}_j . In practice, we can use the empirical distribution of the input feature x_i for μ_i or the uniform distribution. With the sum-to-zero in (1), the functional ANOVA 202 model becomes identifiable, as can be seen in proposition 3.1. Let $\mu = \prod_{i} \mu_{i}$. 203

Proposition 3.1. (Hooker, 2007) Consider two component sets $\{f_S^1, S \in \mathbf{S}\}$ and $\{f_S^2, S \in \mathbf{S}\}$ which satisfy (1). Then, $\sum_{S \in \mathbf{S}} f_S^1(\cdot) \equiv \sum_{S \in \mathbf{S}} f_S^2(\cdot)$ almost everywhere (with respect to μ) if and only if 204 205 $f_S^1(\cdot) \equiv f_S^2(\cdot)$ almost everywhere (with respect to μ) for every $S \in \mathbf{S}$. 206

207 The sum-to-zero condition is not a unique identifiability condition. However, Herren & Hahn (2022) 208 demonstrated that there is an interesting relation between the sum-to-zero condition and SHAP (Lundberg, 2017) which is a well known interpretable AI method. That is, SHAP value of a given 209 input can be calculated easily from the prediction values of each component under regularity con-210 ditions. For a given model f and input vector x, SHAP value of the *j*th input variable is defined 211 as 212

$$\phi_j(f, \mathbf{x}) = \sum_{S \subseteq [p] \setminus \{j\}} \frac{|S|!(p - |S| - 1)!}{p!} (v_f(S \cup \{j\}) - v_f(S)),$$

215

where $v_f(S) = \mathbb{E}[f(\mathbf{X})|\mathbf{X}_S = \mathbf{x}_S]$, where $\mathbf{X} \sim \mu$.

216 **Proposition 3.2.** (Herren & Hahn, 2022) For a given f which is the $GA^{d}M$ satisfying the sum-to-217 zero condition. Then, we have 218

$$\phi_j(f, \mathbf{x}) = \sum_{S \subseteq [p], |S| \le d, j \in S} f_S(\mathbf{x}_S) / |S|.$$
⁽²⁾

The result in equation (2) provides an interesting implication - the functional ANOVA model satisfying the sum-to-zero condition also decomposes SHAP value. That is, the contribution of the 223 interaction between x_i and \mathbf{x}_S to SHAP value $\phi_i(f, \mathbf{x})$ is $f(\mathbf{x}_{S'})/|S'|$, where $S' = S \cup \{j\}$. Note 224 that this interesting relation is not generally valid for identifiability conditions other than the sum-225 to-zero condition.

226 Therefore, for general $GA^{d}M$, we can calculate SHAP value using equation (2), which we refer 227 to as ANOVA-SHAP. One advantage of ANOVA-SHAP is that it is significantly faster to compute 228 compared to Deep-SHAP and Kernel-SHAP proposed by Lundberg (2017), as well as Tree-SHAP 229 proposed by Lundberg et al. (2018). The ANOVA-SHAP experiment results are in Appendix K. 230

231 3.2 ANOVA-NODE 232

219 220 221

222

240 241 242

249 250

251

256 257

259

260 261

262

263 264 265

Incorporating the sum-to-zero condition into existing neural functional ANOVA models would be 233 difficult because standard gradient descent algorithms cannot be applied due to identifiability con-234 straints. The aim of this subsection is to propose a special neural network function for the functional 235 ANOVA model that automatically satisfies the sum-to-zero condition and thus each component can 236 be estimated uniquely by using a standard gradient descent based optimization algorithm. 237

The main idea of the proposed neural network is to model each component as the sum of special but 238 simple neural networks that satisfy the sum-to-zero condition. That is, we set 239

$$f_S(\mathbf{x}_S) = \sum_{k=1}^{K_S} h^S(\mathbf{x}_S | \phi_{S,k}),$$

243 where $\{h^{S}(\cdot | \phi_{S,k})\}_{k=1}^{K_{S}}$ are neural networks with learnable parameters $\phi_{S,k}$ satisfying the sum-to-244 zero condition. 245

ANOVA-NODT. For $h^{S}(\cdot | \phi_{S,k})$, we propose ANOVA-NODT, which is a modification of ODT, 246 as follows. First, we set the depth of h^S to be equal to |S|. For $t \in \{1, ..., |S|\}$, we use the feature 247 function as 248

$$F^t(\mathbf{x}_S) = (\mathbf{x}_S)_t$$

where $(\mathbf{x}_S)_t$ is the *t*th feature of \mathbf{x}_S . To approximate the indicator function I in ODT as a differentiable function, we use $entmax_{\nu}$ as Peters et al. (2019) does. For $b \in \mathbb{R}$ and $\gamma \in \mathbb{R}^+$, we define

$$c_1(x|b,\gamma) := 1 - entmax_{\nu} \left(\left(\frac{x-b}{\gamma}, 0 \right)' \right)_1, \quad c_2(x|b,\gamma) := 1 - c_1(x|b,\gamma),$$

where ν is a hyper-parameter. Finally, we define the ANOVA-NODT model h^S as

$$h^{S}(\mathbf{x}_{S}) = (\mathcal{B}^{S})' \left(\begin{bmatrix} c_{1}((\mathbf{x}_{S})_{1}|b_{1},\gamma_{1}) \\ c_{2}((\mathbf{x}_{S})_{1}|b_{1},\gamma_{1}) \end{bmatrix} \otimes \begin{bmatrix} c_{1}((\mathbf{x}_{S})_{2}|b_{2},\gamma_{2}) \\ c_{2}((\mathbf{x}_{S})_{2}|b_{2},\gamma_{2}) \end{bmatrix} \otimes \cdots \otimes \begin{bmatrix} c_{1}((\mathbf{x}_{S})_{|S|}|b_{|S|},\gamma_{|S|}) \\ c_{2}((\mathbf{x}_{S})_{|S|}|b_{|S|},\gamma_{|S|}) \end{bmatrix} \right)$$

where $\mathcal{B}^{S} \in \mathbb{R}^{2^{|S|}}$, $b_1, \ldots, b_{|S|}$ and $\gamma_1, \ldots, \gamma_{|S|}$ are learnable parameters.

A novel part of ANOVA-NODT is to parameterize \mathcal{B}^S to make h^S always satisfy the sum-to-zero condition. Here, we address the case where |S| = 1, and the case where |S| = 2 is described in Appendix B.2. For the case of $S = \{j\}$, we consider

$$h^{j}(x_{j}) = \beta_{1,j}c_{1}(x_{j}|b_{1,j},\gamma_{1,j}) + \beta_{2,j}c_{2}(x_{j}|b_{1,j},\gamma_{1,j}).$$

266 For a given $\theta_j \in \mathbb{R}$, we let $\beta_{1,j} = \theta_j$ and $\beta_{2,j} = -\frac{\mathbb{E}[c_1(X_j|b_{1,j},\gamma_{1,j})]}{\mathbb{E}[c_2(X_j|b_{1,j},\gamma_{1,j})]}\theta_j$. Then, it is easy to see that 267 $\mathbb{E}[h^j(X_j)] = 0$, meaning that it satisfies the sum-to-zero condition. Hence, we can parameterize 268 $h^{j}(x_{i})$ by $\phi_{i} = (\theta_{i}, b_{1,i}, \gamma_{i})$. Note that there is no constraint on ϕ_{i} and thus standard gradient 269 descent algorithms can be used to learn ϕ_i .



Figure 2: ANOVA-N¹ODE architecture

Finally, for general $S \subseteq [p]$, ANOVA-NODT can be parameterized by

$$\phi_S = \{\theta_S, (b_{1,S}, \gamma_{1,S}), \dots, (b_{|S|,S}, \gamma_{|S|,S})\}$$

287

as

284

285

289 290 291

292 293

294

295 296 297

$$h^{S}(\mathbf{x}_{S}|\phi_{S}) = \sum_{\mathbf{s}\in\{1,2\}^{|S|}} \beta_{\mathbf{s},S}(\phi_{S}) \prod_{i=1}^{|S|} c_{s_{i}}((\mathbf{x}_{S})_{i}|b_{i,S},\gamma_{i,S}),$$

where $\mathbf{s} = (s_1, \dots, s_{|S|})'$ and $\beta_{\mathbf{s},S}(\phi_S)$ are the functions of ϕ_S with $\beta_{\mathbf{1},S}(\phi_S) = \theta_S$.

ANOVA-N^d**ODE.** ANOVA-N^dODE is an ensemble of ANOVA-NODT for the order of interactions up to d. That is, we set

$$f_{ ext{ANOVA-N}^d ext{ODE}}(\mathbf{x}) = \sum_{S \subseteq [p], |S| \le d} \sum_{k=1}^{K_S} h^S(\mathbf{x}_S | \phi_{S,k}),$$

where K_S is the number of ANOVA-NODTs corresponding to the component S. In practice, d = 1, which corresponds to the GAM, and d = 2 are commonly used. The architecture of ANOVA-N¹ODE is shown in Figure 2. Note that ANOVA-N^dODE always satisfies the sum-to-zero condition since each ANOVA-NODT does so. Therefore, standard gradient descent algorithms can be applied without modification. Unless there is any confusion, we write ANOVA-NODE instead of ANOVA-N^dODE for general d.

An interesting theoretical property of ANOVA-NODE is the universal approximation property as the standard neural network has (Hornik et al., 1989). That is, ANOVA-NODE can approximate any arbitrary $GA^{d}M$ function to a desired level of accuracy, as stated in the following theorem.

Theorem 3.3. Let $g_0(\mathbf{x}) := \sum_{S \subset [p], |S| \le d} g_{0,S}(\mathbf{x}_S)$ be a given GA^dM function satisfying the sumto-zero condition. Consider the entmax_{\nu} for $\nu = 1$, and let μ be a measure for any distribution which has a bounded density. If each $g_{0,S}$ is L-Lipschitz continuous¹ for some L > 0, then for any $\epsilon > 0$, there exists $f_{ANOVA-N^d ODE}$ with sufficiently large K_Ss such that

316

323

 $\left\|g_0(\cdot) - f_{ANOVA-N^d ODE}(\cdot)\right\|_{\infty} < \epsilon.$

Training For a given train data and loss function, ANOVA-NODE is trained using any gradient descent algorithm to minimize the empirical risk.

Data preprocessing. To satisfy the sum-to-zero condition in ANOVA-NODT, one element in \mathcal{B}^S is a free parameter, while the remaining elements are parameterized by the product of constants (e.g., $-\mathbb{E}[c_1(X_j|b_{1,j},\gamma_{1,j})]/\mathbb{E}[c_2(X_j|b_{1,j},\gamma_{1,j})])$ and the free parameter. In this case, since the constants can become close to or equal to zero depending on the trained $\{(b_{i,S},\gamma_{1,S}), i = 1, ..., |S|\}$, we scale the data by using a transformation based on the quantiles of a uniform distribution to ensure stable learning.

¹A given function v defined on \mathcal{Z} is L-Lipschitz continuous if $|v(z_1) - v(z_2)| \leq L ||z_1 - z_2||$ for all $z_1, z_2 \in \mathcal{Z}$, where $|| \cdot ||$ is a certain norm defined on \mathcal{Z} .



339 Figure 3: Scatter plots of the stability scores on CALHOUSING dataset. Figures (a) and (d) are the scatter plots for the stability scores of the main effects, where the x-axis is the stability score of ANOVA-N¹ODE and 340 the y-axis is the stability scores of $NA^{1}M$ in (a) and $NB^{1}M$ in (d). Figures (b) and (e) are the scatter plots 341 for the stability scores of the main effects, where the x-axis is the stability score of ANOVA-N²ODE, and the 342 y-axis is the stability score of NA^2M in (b) and NB^2M in (c). Figures (c) and (f) compare the stability scores of 343 the second order interactions of ANOVA-N²ODE to those of NA²M and NB²M. Each dot in the scatter plots 344 corresponds to each component. 345

4 EXPERIMENTS

This section presents the results of numerical experiments. More results along with details about data, algorithms and selection of hyper-parameters are provided in Appendices C to M.

348 349 350

346

347

4.1 STABILITY IN COMPONENT ESTIMATION

ż

351 Similarly to NAM (Agarwal et al., 2021) and NBM (Radenovic et al., 2022), ANOVA-NODE pro-352 vides interpretation through the estimated components. Thus, if the components are not estimated 353 stably, the interpretations based on the estimated components would not be reliable. In this subsection, we investigate the stability of the component estimation of ANOVA-NODE compared with the 354 other baseline models including NAM and NBM. For this purpose, we generate randomly sampled 355 training data and estimate the components of the functional ANOVA model. We repeat this proce-356 dure 10 times to obtain 10 estimates of each component, and measure how similar these 10 estimates 357 are. For the similarity measure, we use

358

359 360

$$\mathcal{SC}(f_S) = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{10} (f_S^j(\mathbf{x}_i) - \bar{f}_S(\mathbf{x}_i))^2}{\sum_{j=1}^{10} (f_S^j(\mathbf{x}_i))^2}$$

361

362 for given pre-selected n many input vectors \mathbf{x}_i , i = 1, ..., n, where f_S^j , j = 1, ..., 10 are the 10 estimates of f_S and \overline{f}_S is their average. A smaller value of $\mathcal{SC}(f_S)$ means a more stable estimation 364 (and thus more stable interpretation). 365

In Figure 3, we present the scatter plots of the stability scores for each estimated component between 366 (NA¹M, NB¹M) vs ANOVA-N¹ODE and (NA²M, NB²M) vs ANOVA-N²ODE on CALHOUSING 367 dataset. It is obvious that our models are more stable than the baselines in component estimation. 368 Consistent results are also observed with other datasets, which are presented in Appendix F.3. Also, 369 the plots of the functional relation of the main effects are provided in Appendix F.1, from which we 370 can feel how much NAM and NBM are unstable in estimating the components. 371

In addition, we compare the overall stability score $SC(f) = \sum_{S} SC(f_S)$. For each of nine bench-372 mark datasets, we calculate the ratios of the overall stability scores of ANOVA-NODE, NAM, and 373 NBM normalized by the overall stability score of ANOVA-NODE, whose results are given in Table 374 1. The unnormalized results of stability score are in Appendix G.1. The results again confirm that 375 ANOVA-NODE is superior in terms of the stability of component estimation. 376

The stability of ANOVA-NODE with respect to the choice of initial values are illustrated in Ap-377 pendix D.2.

Table 2: **Performance of component selection.** We report the averages (standard deviations) of AUROCs of the estimated importance scores of each component on $f^{(1)}$, $f^{(2)}$, $f^{(3)}$ synthetic datasets. The bold faces highlight the best results.

	True model	$f^{(1)}$			$f^{(2)}$			$f^{(3)}$		
-	Models	ANOVA N ² ODE	NA ² M	NB^2M	ANOVA N ² ODE	NA ² M	NB^2M	ANOVA N ² ODE	NA^2M	NB ² M
-	AUROC \uparrow	1.000 (0.00)	0.330 (0.08)	0.522 (0.16)	0.943 (0.01)	0.311 (0.08)	0.481 (0.09)	0.956 (0.02)	0.381 (0.13)	0.477 (0.07)

Table 1: **Stability scores on real datasets.** For each dataset, stability scores of GAM (GA²M) models are normalized by the that of ANOVA-N¹ODE (ANOVA-N²ODE). Lower stability score means more stable interpretation. The bold faces highlight the best results.

		GAM			GA ² M	
Dataset	ANOVA N ¹ ODE	NA^1M	NB^1M	ANOVA N ² ODE	NA^2M	NB^2M
CALHOUSING (Pedregosa et al., 2011a)	1.000	3.750	3.250	1.000	2.209	2.143
WINE(Cortez et al., 2009)	1.000	5.273	3.909	1.000	1.776	1.327
ONLINE (Fernandes et al., 2015)	1.000	2.452	1.742	1.000	1.385	1.385
ABALONE (Warwick et al., 1995)	1.000	1.625	3.250	1.000	1.679	1.357
FICO (fic, 2018)	1.000	1.314	1.314	1.000	1.854	1.563
CHURN(chu, 2017)	1.000	1.588	2.765	1.000	1.894	1.702
CREDIT (cre, 2015)	1.000	3.286	1.190	1.000	2.472	1.472
LETTER (Slate, 1991)	1.000	1.294	0.824	1.000	2.885	1.962
DRYBEAN(dry, 2020)	1.000	2.643	2.500	1.000	1.660	1.528

381 382

384

385

386

4.2 Performance in component selection

An important implication of stable estimation of the components is the ability of selecting signal components. That is, ANOVA-NODE can effectively identify signal components in the true function by measuring the variations of the estimated components. For example, we can consider the l_1 norm of each estimated component (i.e, $||f_S(\mathbf{x}_S)||_1$) as the important score, and select the components whose important scores are large. This simple component selection procedure would not perform well if component estimation is unstable.

405 To investigate how well ANOVA-NODE selects the true signal components, we conduct an experi-406 ment similar to the one in Tsang et al. (2017). We generate synthetic datasets from $Y = f(\mathbf{x}) + \epsilon$, 407 where f is the true prediction model and ϵ is a noise generated from a Gaussian distribution with 408 mean 0 and variance σ_{ϵ}^2 . Then, we apply ANOVA-N²ODE, NA²M and NB²M to calculate the im-409 portance scores of the main effects and second order interactions and examine how well they predict whether a given component is signal. For the performance measure of component selection, we use 410 AUROC obtained from the pairs of $||\hat{f}_S||_1$ and r_S for all $S \subset [p]$ with $|S| \leq 2$, where \hat{f}_S are the estimates of f_S in f and $r_S = \mathbb{I}(||f_S^{(k)}||_1 > 0)$ are the indicators whether f_S are signal or not. 411 412 413

For the true prediction model, we consider the three functions $f^{(k)}$, k = 1, 2, 3 whose details are given in Appendix C.1. We set the data size to 15,000 and set σ_{ϵ}^2 to make the signal-to-noise ratio become 5. Table 2 compares the AUROCs of ANOVA-N²ODE, NA²M and NB²M, which clearly indicates that ANOVA-N²ODE outperforms the baseline models in component selection. More details of component selection with ANOVA-N²ODE are given in Appendix H.

419

4.3 PREDICTION PERFORMANCE

We compare prediction performance of ANOVA-NODE with baseline models. We randomly split
the train, validation and test data into the ratio 70/10/20, where the validation data is used to select
the optimal epoch and the test data is used to measure the prediction performance of the estimated
models. We repeat this random split 10 times to obtain 10 performance measures for prediction. For
the performance measure, we use the Root Mean Square Error (RMSE) for regression datasets and
the Area Under the ROC curve (AUROC) for classification datasets.

Table 3 presents the results of prediction performance of ANOVA-NODE, NODE-GAM, NAM, and
 NBM as well as two black box models. It is obvious that ANOVA-NODE favorably competes with
 its competitors in view of prediction performance. In addition, at the final line, the average ranks of
 each model over the nine datasets are given, which shows that ANOVA-NO²DE exhibits comparable
 prediction performance compared to the baseline models. Details about the experiments are given in
 Appendix C.2.

-				GAN	M			GA ² 1	М		Black l	oox
_	Dataset	Measure	ANOVA N ¹ ODE	NODE GA ¹ M	$\mathrm{NA}^1\mathrm{M}$	NB^1M	ANOVA N ² ODE	NODE GA ² M	NA^2M	NB^2M	XGB	NODE
	CALHOUSING	RMSE	0.614	0.581	0.659	0.594	0.512	0.515	0.525	0.502	0.452	0.482
	CALIFOUSING	KINDL 4	(0.01)	(0.01)	(0.01)	(0.08)	(0.01)	(0.01)	(0.02)	(0.03)	(0.01)	(0.01)
	WINE	RMSE	0.725	0.723	0.733	0.724	0.704	0.730	0.720	0.702	0.635	0.646
	WINE	KINDL V	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.03)	(0.03)	(0.03)
	ONLINE	RMSE	1.111	1.121	1.350	1.187	1.111	1.137	1.313	1.179	1.122	1.112
	ONLINE	KINDL V	(0.25)	(0.27)	(0.57)	(0.25)	(0.25)	(0.26)	(0.46)	(0.21)	(0.26)	(0.27)
	A BALONE	RMSE	2.135	2.141	2.171	2.167	2.087	2.100	2.088	2.088	2.164	2.086
_	ABALONE	KNDL 4	(0.09)	(0.09)	(0.08)	(0.09)	(0.08)	(0.10)	(0.08)	(0.08)	(0.09)	(0.09)
	FICO		0.799	0.795	0.788	0.797	0.800	0.793	0.799	0.799	0.796	0.795
	FICO	AUROC T	(0.007)	(0.009)	(0.006)	(0.006)	(0.007)	(0.007)	(0.007)	(0.008)	(0.008)	(0.008)
	CHURN		0.839	0.824	0.846	0.845	0.842	0.830	0.844	0.844	0.846	0.844
	CHUKN	AUROC	(0.012)	(0.012)	(0.011)	(0.012)	(0.012)	(0.011)	(0.011)	(0.011)	(0.012)	(0.013)
	CREDIT		0.983	0.983	0.976	0.972	0.984	0.985	0.980	0.985	0.983	0.984
	CREDIT	AUROC	(0.005)	(0.005)	(0.012)	(0.011)	(0.006)	(0.006)	(0.007)	(0.004)	(0.004)	(0.009)
	LETTER		0.900	0.910	0.904	0.910	0.984	0.988	0.986	0.990	0.997	0.998
	LEITER	AUKOC	(0.003)	(0.002)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
	DRVDEAN		0.995	0.996	0.996	0.994	0.998	0.996	0.995	0.995	0.997	0.996
_	DKIBEAN	AUKOC T	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
		Rank avg \downarrow	6.33	5.56	8.00	7.44	3.33	5.67	5.44	3.67	3.44	2.89

Table 3: Prediction performance. We report the averages (standard deviations) of the prediction performance
 measure. In addition, we report the averages of ranks of prediction performance of each model on nine datasets.
 The optimal (or suboptimal) results are highlighted in **bold** (or underlined).

4.4 APPLICATION TO HIGH-DIMENSIONAL DATA

452 To see whether ANOVA-NODE is applicable to high-dimensional data, we analyze three additional 453 datasets with input dimensions ranging from 136 to 699. See Table 8 of Appendix for details of 454 these three datasets. For ANOVA-N¹ODE, we include all main effects into the model. For ANOVA-455 N^2 ODE, however, the number of second order interactions is too large so that considering all the 456 main effects and second order interactions would be difficult unless very large computing resources 457 are available. A simple alternative is to screen out unnecessary second order interactions a priori and include only selected second order interactions (and all the main effects) into the model. In 458 the experiment, we use Neural Interaction Detection (NID, Tsang et al. (2017)) for the interaction 459 screening. The number of selected interactions is given in Appendix C.2. 460

From Table 4 and Table 5, we observe that ANOVA-NODE shows favorable prediction performance
compared with NAM and NBM and estimates the components more stably on high-dimensional
datasets. In addition, note that the RMSE of NB²M with all second order interactions on MICROSOFT
dataset is reported as 0.750 by Radenovic et al. (2022). That is, screening interactions using NID
does not hamper prediction performance much.

Table 4: **Prediction performance on high-dimensional datasets.** We report the averages (standard deviations) of the prediction performance for 10 randomly sampled training data from the high-dimensional datasets. The bold faces highlight the best results.

			GAM			GA ² M	
Dataset	Measure	ANOVA N ¹ ODE	NA^1M	NB^1M	$NID + \frac{ANOVA}{N^2ODE}$	$NID + NA^2M$	$NID + NB^2M$
MICROSOFT	$RMSE \downarrow$	0.756 (0.001)	0.774 (0.001)	0.770 (0.001)	0.754 (0.001)	0.761 (0.001)	0.755 (0.001)
Yahoo	$RMSE \downarrow$	0.787 (0.002)	0.797 (0.002)	0.783 (0.002)	0.779 (0.001)	0.793 (0.002)	0.779 (0.002)
MADELON	AUROC ↑	0.587 (0.02)	0.587 (0.02)	0.582 (0.03)	0.605 (0.01)	0.568 (0.03)	0.594 (0.02)

473 474 475

476

477

478

479

480

481

482

466

467

468

469 470

471 472

4.5 ANOVA-NODE WITH MONOTONE CONSTRAINTS

Monotone constraint. In practice, prior knowledge that some main effects are monotone functions are available and it is needed to reflect this prior knowledge in the training phase. A notable example is the credit scoring model where certain input features should have monotone main effects (Chen & Li, 2014; Chen & Ye, 2022). An additional advantage of ANOVA-NODE is to accommodate the monotone constraints in the model easily. Suppose that f_i is monotonically increasing.

Then, ANOVA-NODE can estimate f_j monotonically increasingly by letting the θ_j in ANOVA-NODT $h_j(x_j|\phi_j)$ be less than or equal to 0. See Appendix B.1 for details.

483 484

Application to Image data. Monotone constraint helps avoiding unreasonable interpretation. To illustrate this advantage, we conduct an experiment with an image dataset. We use CELEBA (Liu

Table 5: **Stability scores on the high-dimensional datasets.** For each dataset, stability scores of of GAM (GA²M) models normalized by the that of ANOVA-N¹ODE (ANOVA-N²ODE) are presented. Lower stability scores imply more stable interpretation. The bold faces highlight the best results.



Figure 4: Plots of the functional relations of 'No Beard' and 'Wearing Lipstick' on CELEBA dataset estimated by ANOVA-N¹ODE with and without the monotone constraint.

et al., 2015) dataset which has 40 binary attributes for each image. To apply ANOVA-NODE to
CELEBA dataset, we consider Concept Bottleneck Model (CBM, Koh et al. (2020)) similar to the
one used in Radenovic et al. (2022). In CBM, rather than directly inputting the embedding vector
derived from an image data through a CNN into a classifier, the CNN initially predicts each concept
accompanied with each image. Then, these predicted values of each concept are subsequently used
as the input of a DNN classifier. For our experiment, we use a pretrained ResNet18, where the last
layer consists of a linear transformation with a softmax activation function, and we replace the final
DNN classifier with ANOVA-NODE.

Among the attributes, we set 'gender' as the target label and the remaining attributes are set as concepts for images. Since 'male' is labeled as 1 and 'female' as 0, a higher value of each component results in a higher chance of being classified as 'male'.

Figure 4 presents two functional relations of the main effects of the concepts 'No Beard' and 'Wearing Lipstick', estimated on a randomly sampled training dataset with and without the monotone constraint. Note that the functional relations are quite different even though their prediction performances, which is given in Table 15 of Appendix E.1.2, are similar. It is a common sense that an image having the concept of 'No Beard' and 'Wearing Lipstick' has a higher chance of being a female and thus the functional relations are expected to be decrease. Figure 4 illustrates that a completely opposite result to our common sense could be obtained in practice. Implications of the opposite functional relations to interpretation of each image are discusses in Appendix E.1.2.

522 523 524

489

490

491

492

493

494

495

496

497

498

499

500

501

502

4.6 ADDITIONAL EXPERIMENTS

In Appendix L, we confirm that the component estimation of ANOVA-N²ODE becomes highly
 unstable when the sum-to-zero condition is not imposed, and in Appendix M, we discuss a method
 for enforcing the sum-to-zero condition after training NAM or NBM.

527 528 529

530

5 CONCLUSION

In this paper, we propose a novel XAI model called ANOVA-NODE for estimating the functional
 ANOVA model stably based on Neural Oblivious Decision Tree (NODT). We theoretically demon strate that ANOVA-NODE can approximate a smooth function well. We also empirically show that
 prediction performance of ANOVA-NODE is comparable to its competitors.

One way to make ANOVA-NODE computationally more efficient is to combine ANOVA-NODE
and NBM (Radenovic et al., 2022), which has significantly fewer weight parameters. We explored
this approach, and the algorithm and results of the experiment are given in Appendix G.2. Even
though this algorithm is computationally more efficient than ANOVA-NODE itself, analyzing highdimensional data is still computationally demanding due to too many components. In general,
scaling-up interpretable AI algorithms is a promising future research topic.

Reproducibility Statement. The complete proofs of the theoretical results are presented thor oughly and rigorously in Appendix A. Details regarding the experimental implementation, datasets,
 libraries, and hyper-parameters are outlined in Appendix C. Furthermore, the proposed ANOVA NODE in this paper is an ensemble of ANOVA-NODT, and Appendix B provides a detailed expla nation of ANOVA-NODT. Therefore, by referring to Appendix B, ANOVA-NODT can be imple mented for each component S, and the ensemble ANOVA-NODE can also be easily implemented.
 Finally, We will provide the source code and make it open access by uploading it on the web after
 acceptance.

549 REFERENCES

548

550

563

564

- 551 Credit card fraud detection. kaggle, 2015. https://www.kaggle.com/datasets/mlg 552 ulb/creditcardfraud.
- Telco customer churn. kaggle, 2017. https://www.kaggle.com/datasets/blastchar/telco-customerchurn/data.
- 556 Fico heloc. FICO Explainable Learning Challenge, 2018. https://community.fico.com/s/ 557 explainable-machine-learning-challenge.
- ⁵⁵⁸ Dry bean. UCI Machine Learning Repository, 2020. DOI: https://doi.org/10.24432/C50S4B.
- Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich Caruana, and Geoffrey E Hinton. Neural additive models: Interpretable machine learning with neural nets.
 Advances in neural information processing systems, 34:4699–4711, 2021.
 - Mathieu Blondel, Akinori Fujino, Naonori Ueda, and Masakazu Ishihata. Higher-order factorization machines. *Advances in neural information processing systems*, 29, 2016.
- Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. Node-gam: Neural generalized additive
 model for interpretable deep learning. *arXiv preprint arXiv:2106.01613*, 2021.
- Gaëlle Chastaing, Fabrice Gamboa, and Clémentine Prieur. Generalized hoeffding-sobol decompo sition for dependent variables-application to sensitivity analysis. 2012.
- 573 Chih-Chuan Chen and Sheng-Tun Li. Credit rating with a monotonicity-constrained support vector
 574 machine model. *Expert Systems with Applications*, 41(16):7235–7247, 2014.
- Dangxing Chen and Weicheng Ye. Monotonic neural additive models: Pursuing regulated machine learning models for credit scoring. In *Proceedings of the third ACM international conference on AI in finance*, pp. 70–78, 2022.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794,
 2016.
- Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. 2010.
- Paulo Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Wine Quality. 2009. DOI: https://doi.org/10.24432/C56S3T.
- 587
 588
 589
 Kelwin Fernandes, Pedro Vinagre, Paulo Cortez, and Pedro Sernadela. Online News Popularity. UCI Machine Learning Repository, 2015. DOI: https://doi.org/10.24432/C5NS3V.
- Chong Gu and Grace Wahba. Smoothing spline anova with component-wise bayesian "confidence intervals". *Journal of Computational and Graphical Statistics*, 2(1):97–117, 1993.
- 593 Isabelle Guyon. Madelon. UCI Machine Learning Repository, 2004. DOI: https://doi.org/10.24432/C5602H.

- 594 Trevor Hastie and Robert Tibshirani. Generalized additive models: some applications. Journal of the American Statistical Association, 82(398):371–386, 1987. 596 Andrew Herren and P Richard Hahn. Statistical aspects of shap: Functional anova for model inter-597 pretation. arXiv preprint arXiv:2208.09970, 2022. 598 Wassily Hoeffding. A class of statistics with asymptotically normal distribution. Breakthroughs in 600 statistics: Foundations and basic theory, pp. 308–334, 1992. 601 Giles Hooker. Generalized functional anova diagnostics for high-dimensional functions of depen-602 dent variables. Journal of computational and graphical statistics, 16(3):709-732, 2007. 603 604 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are uni-605 versal approximators. *Neural networks*, 2(5):359–366, 1989. 606 Jianhua Z Huang, Charles Kooperberg, Charles J Stone, and Young K Truong. Functional anova 607 modeling for proportional hazards regression. Annals of Statistics, pp. 961–999, 2000. 608 609 Jinseog Kim, Yongdai Kim, Yuwon Kim, Sunghoon Kwon, and Sangin Lee. Boosting on the func-610 tional anova decomposition. Statistics and Its Interface, 2(3):361–368, 2009. 611 Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and 612 Percy Liang. Concept bottleneck models. In International conference on machine learning, pp. 613 5338-5348. PMLR, 2020. 614 615 Ron Kohavi and Chia-Hsin Li. Oblivious decision trees, graphs, and top-down pruning. In IJCAI, 616 volume 2, pp. 1071–1077, 1995. 617 Benjamin Lengerich, Sarah Tan, Chun-Hao Chang, Giles Hooker, and Rich Caruana. Purifying 618 interaction effects with the functional anova: An efficient algorithm for recovering identifiable 619 additive models. In International Conference on Artificial Intelligence and Statistics, pp. 2402– 620 2412. PMLR, 2020. 621 Yi Lin and Hao Helen Zhang. Component selection and smoothing in multivariate nonparametric 622 regression. The Annals of statistics, 34(5):2272-2297, 2006. 623 624 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. 625 In Proceedings of the IEEE international conference on computer vision, pp. 3730–3738, 2015. 626 Scott Lundberg. A unified approach to interpreting model predictions. arXiv preprint 627 arXiv:1705.07874, 2017. 628 SM Lundberg, GG Erion, and SI Lee. Consistent individualized feature attribution for tree ensem-630 bles. arxiv 2018. arXiv preprint arXiv:1802.03888, 10, 2018. 631 Kaspar Märtens and Christopher Yau. Neural decomposition: Functional anova with variational 632 autoencoders. In International Conference on Artificial Intelligence and Statistics, pp. 2917-633 2927. PMLR, 2020. 634 635 Christoph Molnar. Interpretable machine learning. Lulu. com, 2020. 636 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-637 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and 638 E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 639 12:2825-2830, 2011a. 640 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier 641 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: 642 Machine learning in python. the Journal of machine Learning research, 12:2825–2830, 2011b. 643 644 Ben Peters, Vlad Niculae, and André FT Martins. Sparse sequence-to-sequence models. arXiv 645 preprint arXiv:1905.05702, 2019. 646
- 647 Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for deep learning on tabular data. arxiv 2019. *arXiv preprint arXiv:1909.06312*, 2019.

648 649	Tao Qin and Tie-Yan Liu. Introducing LETOR 4.0 datasets. CoRR, abs/1306.2597, 2013. URL http://arxiv.org/abs/1306.2597.
650 651 652	Filip Radenovic, Abhimanyu Dubey, and Dhruv Mahajan. Neural basis models for interpretability. Advances in Neural Information Processing Systems, 35:8414–8426, 2022.
653 654 655	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pp. 1135–1144, 2016.
657 658	David Rügamer. Scalable higher-order tensor product spline models. In International Conference on Artificial Intelligence and Statistics, pp. 1–9. PMLR, 2024.
659 660 661	Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In <i>Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society</i> , pp. 180–186, 2020.
662 663 664	David Slate. Letter Recognition. UCI Machine Learning Repository, 1991. DOI: https://doi.org/10.24432/C5ZP40.
665 666	Michael Tsang, Dehua Cheng, and Yan Liu. Detecting statistical interactions from neural network weights. <i>arXiv preprint arXiv:1705.04977</i> , 2017.
667 668 669	Nash Warwick, Sellers Tracy, Talbot Simon, Cawthorn Andrew, and Ford Wes. Abalone. UCI Machine Learning Repository, 1995. DOI: https://doi.org/10.24432/C55C7W.
670 671 672	Yahoo. Yahoo! learning to rank challenge. <i>https://webscope.sandbox.yahoo.com/catalog.php?datatype=c</i> , 2010.
673 674	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
600	
601	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

Supplementary material

Appendix

A PROOFS FOR THEORETICAL RESULTS

A.1 PROOF OF PROPOSITION 3.1

For a given function f, consider two component functions sets $\{f_S^1, S \in \mathbf{S}\}$ and $\{f_S^2, S \in \mathbf{S}\}$ such that

$$f(\mathbf{x}) = \sum_{S \in \mathbf{S}} f_S^1(\mathbf{x}_S) = \sum_{S \in \mathbf{S}} f_S^2(\mathbf{x}_S)$$

for every $\mathbf{x} \in \mathcal{X}$, where each component function satisfies the sum-to-zero condition (1).

For any $S, V \in \mathbf{S}$ such that $S \neq V$, we have $S \setminus V \neq \emptyset$ or $V \setminus S \neq \emptyset$. Assume $S \setminus V \neq \emptyset$ without loss of generality. For $i_1, i_2 \in \{1, 2\}$ and $s \in S \setminus V$, we get

$$\int_{\mathcal{X}} f_{S}^{i_{1}}(\mathbf{X}_{S}) f_{V}^{i_{2}}(\mathbf{X}_{V}) d\Pi_{j=1}^{p} \mu_{j}$$

$$= \int_{\mathcal{X}_{[p] \setminus \{s\}}} \left[\int_{\mathcal{X}_{s}} f_{S}^{i_{1}}\left(\mathbf{x}_{S \setminus \{j\}} = \mathbf{X}_{S \setminus \{j\}}, \mathbf{x}_{S} = \mathbf{X}_{S} \right) d\mu_{s} \right] f_{V}^{i_{2}}(\mathbf{X}_{V}) d\Pi_{j \neq s} \mu_{j}$$

$$= 0$$

by the sum-to-zero condition, and hence

$$\int_{\mathcal{X}} (f_S^1(\mathbf{X}_S) - f_S^2(\mathbf{X}_S)) (f_V^1(\mathbf{X}_V) - f_V^2(\mathbf{X}_V)) d\Pi_{j=1}^p \mu_j = 0.$$

Then, we obtain

$$\begin{split} &\sum_{S \in \mathbf{S}} \int_{\mathcal{X}} (f_S^1(\mathbf{X}_S) - f_S^2(\mathbf{X}_S))^2 d\Pi_{j=1}^p \mu_j \\ &= \sum_{S \in \mathbf{S}} \int_{\mathcal{X}} (f_S^1(\mathbf{X}_S) - f_S^2(\mathbf{X}_S))^2 d\Pi_{j=1}^p \mu_j \\ &\quad + \sum_{S,V \in \mathbf{S}, S \neq V} \int_{\mathcal{X}} 2(f_S^1(\mathbf{X}_S) - f_S^2(\mathbf{X}_S))(f_V^1(\mathbf{X}_V) - f_V^2(\mathbf{X}_V)) d\Pi_{j=1}^p \mu_j \\ &= \int_{\mathcal{X}} \left[\sum_{S \in \mathbf{S}} \left(f_S^1(\mathbf{X}_S) - f_S^2(\mathbf{X}_S) \right) \right]^2 d\Pi_{j=1}^p \mu_j \\ &= \int_{\mathcal{X}} \left[f(\mathbf{X}) - f(\mathbf{X}) \right]^2 d\Pi_{j=1}^p \mu_j \\ &= 0. \end{split}$$

To sum up, we have $f_S^1(\cdot) \equiv f_S^2(\cdot)$ almost everywhere for $S \in \mathbf{S}$, which completes the proof.

756 A.2 PROOF OF THEOREM 3.3

758 **Challenging part of Theorem 3.3.** The most important and challenging part of Theorem 3.3 is 759 decomposing the ensemble function $f_{\mathcal{E}}$, which approximates the true function well, into ANOVA-760 NODTs.

Case of d = 1. It is enough to show that for every $j \in [p]$, there exists a set of ANOVA-NODTs $\{h^{S}(\cdot | \phi_{j,k})\}_{k=1}^{K_{S}}$ such that

$$\left\|g_{0,j}(\cdot) - \sum_{k=1}^{K_S} h^S(\cdot|\phi_{j,k})\right\|_{\infty} < \frac{\epsilon}{p}.$$

768 769

761 762

763

We denote $\sigma(x) := 1/(1 + \exp(-x))$ as the sigmoid function. We consider $\nu = 1$ for simplicity, which results in $entmax_{\nu}((x,0)') = \sigma(x)$. Also, we assume that μ_j for $j \in [p]$ admits a density with respect to Lebesgue measure which is bounded above and below. In cases where the density does not exist, such as the empirical distribution, we can construct a *K*-equal-sized partition of \mathcal{X}_j and then handle regions with zero measure and non-zero measure separately, similar to the proof described below.

T76 Let $0 < p_L < p_R < \infty$ be the lower and upper bound of the density of μ_j , respectively. Let $\{\Omega_k\}_{k=1}^K = \{[\chi_{k-1}, \chi_k)\}_{k=1}^K$ be a interval partition of \mathcal{X}_j such that $\mu_j(\Omega_k) = \frac{1}{K}$. We have $|\chi_k - \chi_{k-1}| \le \frac{1}{p_L K}$ for $k \in [K]$. For $\gamma = 1/K^3$, we define $\ell_k(\cdot)$ as

780

781 782

785

786 787

791

792 793 794

807 808

809

788 Note that for every $x \in \mathcal{X}_j$, $\sum_{k=1}^{K} \ell_k(x) = 1$ and $0 \le \ell_k(\cdot) \le 1$ holds for every $k \in [K]$. Also, 789 $\{\ell_k\}_{k=1}^{K}$ have the following properties.

Lemma A.1. For any $k \in [K]$, we have

 $\ell_1(x) = 1 - \sigma\left(\frac{x - \chi_1}{\gamma}\right),$

 $\ell_K(x) = \sigma\bigg(\frac{x - \chi_{K-1}}{\gamma}\bigg).$

 $\ell_k(x) = \sigma\left(\frac{x-\chi_{k-1}}{\gamma}\right) - \sigma\left(\frac{x-\chi_k}{\gamma}\right),$

$$\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \notin \Omega_k)] < \frac{3p_U}{K^2}$$

 $k \in \{2, \dots, K-1\}$

and

$$\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)] > \frac{1}{3K}.$$

The proof of Lemma A.1 is provided in Section A.3. Now, we consider the ensemble function

$$f_{\mathcal{E}}(x) = \sum_{k=1}^{K} \delta_k \ell_k(x)$$

where δ_k is defined as

 $\delta_k = \frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)g_0(X_j)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)]}.$

For any $x \in \mathcal{X}_i$, we have $\left|g_0(x) - f_{\mathcal{E}}(x)\right| = \left|g_0(x) - \sum_{k=1}^{K} \delta_k \ell_k(x)\right|$ $= \left| \sum_{k=1}^{K} (g_0(x) - \delta_k) \ell_k(x) \right|$ $\leq \sum_{k=1}^{K} |g_0(x) - \delta_k| \ell_k(x)$ $=\sum_{k=1}^{K} \left| g_0(x) - \frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)g_0(X_j)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)]} \right| \ell_k(x)$ $=\sum_{k=1}^{K} \left| \frac{\mathbb{E}_{\mathcal{X}_{j}}[\ell_{k}(X_{j})g_{0}(x)] - \mathbb{E}_{\mathcal{X}_{j}}[g_{0}(X_{j})\ell_{k}(X_{j})]}{\mathbb{E}_{\mathcal{X}_{j}}[\ell_{k}(X_{j})]} \right| \ell_{k}(x)$ $\leq L \sum_{k=1}^{K} \left| \frac{\mathbb{E}_{\mathcal{X}_{j}}[\ell_{k}(X_{j})|x - X_{j}|]}{\mathbb{E}_{\mathcal{X}_{j}}[\ell_{k}(X_{j})]} \right| \ell_{k}(x).$ (3)

Let
$$r \in [K]$$
 is the index of partition such that $x \in [\chi_{r-1}, \chi_r)$. For $k \in \{r-1, r, r+1\}$, we have

$$\frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)|x - X_j|]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)]} \leq \frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)|x - X_j|\mathbb{I}(X_j \in \Omega_k)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)]} + \frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)|x - X_j|\mathbb{I}(X_j \notin \Omega_k)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)]}$$

$$\leq \frac{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)(\frac{2}{p_LK})\mathbb{I}(X_j \in \Omega_k)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)]} + \frac{2a \cdot \mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \notin \Omega_k)]}{\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)]}$$

$$\leq \frac{2}{p_LK} + \frac{12aC_{max}}{K}$$

$$\leq \frac{C'}{K},$$

for some constant C' > 0, where the third inequality holds by Lemma A.1. For $k \le r - 2$, we have $x \ge \chi_{r-1}$ and $\chi_k \le \chi_{r-2}$ and hence

Also, for $k \ge r+2$, we have $x \le \chi_r$ and $\chi_{k-1} \ge \chi_{r+1}$ and hence

$$|\ell_k(x)| \le \sigma \left(\frac{x - \chi_{k-1}}{\gamma}\right)$$

$$|\ell_k(x)| \le \sigma \left(\frac{x - \chi_{k-1}}{\gamma}\right)$$

$$\le \sigma \left(\frac{\chi_r - \chi_{r+1}}{\gamma}\right)$$

$$< \frac{1}{1 + \exp(K)}.$$

To sum up, we get

861
862
863

$$|g_0(x) - f_{\mathcal{E}}(x)| \le (3)$$

 $\le L\left(\frac{C'}{K} + \frac{K-1}{1 + \exp(K)}\right)$

Now, for $f_{\mathcal{E}}(x) = \sum_{k=1}^{K} \delta_k \ell_k(x)$, our goal is to show that there exists a set of ANOVA-NODTs $\{h^S(\cdot | \phi_{j,k})\}_{k=1}^{K_S}$ such that

$$\sum_{j=1}^{K_S} h^S(x|\phi_{S,j}) = f_{\mathcal{E}}(x)$$

for every $x \in \mathcal{X}_j$. To derive the result, we use the following lemma.

Lemma A.2. For every $T \in \{2, \ldots, K\}$, we define

$$\ell_{1,T}(x) = 1 - \sigma\left(\frac{x - \chi_1}{\gamma}\right),$$

$$\ell_{t,T}(x) = \sigma\left(\frac{x - \chi_{t-1}}{\gamma}\right) - \sigma\left(\frac{x - \chi_t}{\gamma}\right), \qquad t \in \{2, \dots, T-1\}$$

$$\ell_{T,T}(x) = \sigma\left(\frac{x - \chi_{T-1}}{\gamma}\right).$$

Then, for any given $T \in \{3, \ldots, K\}$ and for any ρ_1, \ldots, ρ_T satisfying

$$\mathbb{E}_{\mathcal{X}_j}\left[\sum_{t=1}^T \rho_t \ell_{t,T}(X_j)\right] = 0,\tag{4}$$

there exist $\kappa_1, \ldots, \kappa_{T-1}$, η and τ such that

$$\sum_{t=1}^{T} \rho_t \ell_{t,T}(x) = \sum_{t=1}^{T-1} \kappa_t \ell_{t,T-1}(x) + \left[\eta \cdot c_1(x|\chi_{T-1},\gamma) + \tau \cdot c_2(x|\chi_{T-1},\gamma)\right],\tag{5}$$

$$\mathbb{E}_{\mathcal{X}_{j}}\left[\sum_{t=1}^{T-1} \kappa_{t}\ell_{t,T-1}(X_{j})\right] = 0, \quad \mathbb{E}_{\mathcal{X}_{j}}\left[\eta \cdot c_{1}(X_{j}|\chi_{T-1},\gamma) + \tau \cdot c_{2}(X_{j}|\chi_{T-1},\gamma)\right] = 0.$$
(6)

The proof of Lemma A.2 is provided in Section A.3. Since

893
894
895
896
897
898
899
900
900
900
901
901
902
903
904
905
905
906
907
907
907
907
907
908

$$\mathbb{E}_{\chi_j}[f_{\mathcal{E}}(X_j)] = \mathbb{E}_{\chi_j}[\ell_k(X_j)g_0(X_j)] = \mathbb{E}_{\chi_j}[\ell_k(X_j)g_0(X_j)]$$

we can decompose $f_{\mathcal{E}}(x)$ using Lemma A.2 by numerical induction. Note that $\rho_1 \ell_{1,2}(\cdot) + \rho_2 \ell_{1,2}(\cdot)$ with $\mathbb{E}_{\mathcal{X}_j}[\rho_1 \ell_{2,2}(X_j) + \rho_2 \ell_{1,2}(X_j)] = 0$ is an ANOVA-NODT and for $k \in \{2, \ldots, K\}$, $\eta \cdot c_1(\cdot|\chi_{T-1}, \gamma) + \tau \cdot c_2(\cdot|\chi_{T-1}, \gamma)$ with $\mathbb{E}_{\mathcal{X}_j}[\eta \cdot c_1(X_j|\chi_{T-1}, \gamma) + \tau \cdot c_2(X_j|\chi_{T-1}, \gamma)]$ is also an ANOVA-NODT. Hence, we can find a set of ANOVA-NODTs $\{h^S(\cdot |\phi_{j,k})\}_{k=1}^{K_S}$ such that

$$\left\| g_{0,j} - \sum_{k=1}^{K_S} h^S(x|\phi_{j,k}) = f_{\mathcal{E}}(x) \right\|_{\infty} < L\left(\frac{C'}{K} + \frac{K-1}{1 + \exp(K)}\right).$$

915 By choosing sufficiently large K, we obtain the assertion.

918 A.3 PROOFS FOR AUXILIARY LEMMAS

Lemma A.1. For $x \leq \chi_{k-1} - \frac{1}{K^2}$, we have

$$|\ell_k(x)| \le \sigma \left(\frac{x - \chi_{k-1}}{\gamma}\right)$$
$$\le \sigma \left(-\frac{1}{K^2 \gamma}\right)$$
$$= \frac{1}{1 + \exp(K)}.$$

Also, for $x \ge \chi_k + \frac{1}{K^2}$, we have

$ \ell_k(x) $	$ \le 1 - \sigma \left(\frac{x - \chi_k}{\gamma} \right)$
	$\leq 1 - \sigma \left(\frac{1}{K^2 \gamma}\right)$
	$=\frac{1}{1+\exp(K)}.$

Hence, we obtain

$$\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \notin \Omega_k)] \leq \mathbb{P}\left(X_j \in \left(\chi_{k-1} - \frac{1}{K^2}, \chi_{k-1}\right) \cup \left(\chi_k, \chi_k + \frac{1}{K^2}\right)\right) + \frac{1}{1 + \exp(K)}$$
$$< \frac{3p_U}{K^2}.$$

Also, for $x \in \left[\chi_{k-1} + \frac{1}{K^2}, \chi_k - \frac{1}{K^2}\right]$, we have

$$\ell_k(x) \ge \sigma\left(\frac{x - \chi_{k-1}}{\gamma}\right) - \sigma\left(\frac{x - \chi_k}{\gamma}\right)$$
$$\ge \sigma(K) - \sigma(-K)$$
$$> \frac{1}{2}$$

for sufficiently large K. Hence, we obtain

$$\mathbb{E}_{\mathcal{X}_j}[\ell_k(X_j)\mathbb{I}(X_j \in \Omega_k)] > \mathbb{P}\left(X_j \in \left[\chi_{k-1} + \frac{1}{K^2}, \chi_k - \frac{1}{K^2}\right]\right) \cdot \frac{1}{2}$$
$$> \frac{1}{3K}.$$

L			L
L			

Lemma A.2. We define

967
968
969
970
971

$$\eta := -\mathbb{E}_{\mathcal{X}_j}[\ell_{T,T}(X_j)](\rho_T - \rho_{T-1}),$$

 $\tau := \sum_{t=1}^{T-1} \mathbb{E}_{\mathcal{X}_j}[\ell_{t,T}(X_j)](\rho_T - \rho_{T-1}),$
 $\kappa_t := \rho_t - \eta,$
 $t \in [T-1].$

 $\sum_{t=1}^{T-1} \kappa_t \ell_{t,T-1}(x) + [\eta \cdot c_1(x|\chi_{T-1},\gamma) + \tau \cdot c_2(x|\chi_{T-1},\gamma)]$

$$= \sum_{t=1}^{T-2} (\rho_t - \eta) \ell_{t,T}(x) + (\rho_{T-1} - \eta) \cdot \sigma \left(\frac{2}{2}\right)$$
$$= \sum_{t=1}^{T-1} (\rho_t - \eta) \ell_{t,T}(x)$$
$$= \sum_{t=1}^{T} (\rho_t - \eta) \ell_{t,T}(x)$$

Then, we obtain (5) by

where the second equality holds by

$$\tau - \eta = \sum_{t=1}^{T} \mathbb{E}_{\mathcal{X}_j}[c_t(X_j)](\rho_T - \rho_{T-1}) = \rho_T - \rho_{T-1}.$$

we obtain (6) by (4) and (5).

 $+(\rho_{T-1}-\eta)\cdot\left(\sigma\left(\frac{x-\chi_{T-2}}{\gamma}\right)-\sigma\left(\frac{x-\chi_{T-1}}{\gamma}\right)\right)$ $\frac{x - \chi_{T-1}}{\gamma} + (\rho_T - \rho_{T-1}) \cdot \sigma \left(\frac{x - \chi_{T-1}}{\gamma}\right) + \eta$ $+(\rho_T-\eta)\cdot\sigma\left(\frac{x-\chi_{T-1}}{\gamma}\right)+\eta$ $+\eta$ $=\sum_{t=1}^{\infty}\rho_t\ell_{t,T}(x),$

 $=\sum_{t=1}^{T-2} (\rho_t - \eta)\ell_{t,T}(x) + (\rho_{T-1} - \eta) \cdot \sigma\left(\frac{x - \chi_{T-2}}{\gamma}\right) + (\tau - \eta) \cdot \sigma\left(\frac{x - \chi_{T-1}}{\gamma}\right) + \eta$

Also, since we have

 $\mathbb{E}_{\mathcal{X}_i}[\eta \cdot c_1(x|\chi_{T-1},\gamma) + \tau \cdot c_2(x|\chi_{T-1},\gamma)]$ $= \eta \cdot \left(1 - \mathbb{E}_{\mathcal{X}_i}[\ell_{T,T}(X_i)]\right) + \tau \cdot \mathbb{E}_{\mathcal{X}_i}[\ell_{T,T}(X_i)]$ $= \eta + \mathbb{E}_{\mathcal{X}_i}[\ell_{T,T}(X_i)](\tau - \eta)$ $= -\mathbb{E}_{\chi_{i}}[\ell_{T,T}(X_{i})](\rho_{T} - \rho_{T-1}) + \mathbb{E}_{\chi_{i}}[\ell_{T,T}(X_{i})](\rho_{T} - \rho_{T-1})$ = 0,

Case of d = 2. The proof for the case of d = 2 is too complex to express by mathematical notation, so we explain the process of decomposing the ensemble function $f_{\mathcal{E}}$ into ANOVA-NODTs, which is the most important part of Theorem 3.3, using a top example. Consider the ensemble function $f_{\mathcal{E}}(x_1, x_2)$ which is defined as below, with 9 height parameters: We explain the process of decomposing $f_{\mathcal{E}}$ into 4 ANOVA-NODTs in 2 steps.

<i>b</i> ₁₂	β ₁₁	β_{12}	β_{13}		$\beta_{11} - \alpha_1$	η_3		α1	
L	β_{12}	β_{22}	β_{23}	=	$eta_{12}-lpha_2$	η_2	+	$lpha_2$	
b_{22}	β_{13}	β_{23}	β_{33}		$eta_{13}-lpha_3$	η_1		α3	
		+ / + / + /	$\beta_{13} \left(1 - \alpha \right)$ $\beta_{23} \left(\sigma \left(\frac{x}{2} \right)$ $\beta_{33} \sigma \left(\frac{x_1}{2} \right)$	$ \frac{\sigma\left(\frac{x_1-b}{\gamma_{11}},\frac{b}{\gamma_{11}},\frac{b}{\gamma_{11}},\frac{b}{\gamma_{11}},\frac{b}{\gamma_{21}}\right) \sigma $	$\left(\frac{x_{1}}{\gamma_{22}}\right) \sigma\left(\frac{x_{1}}{\gamma_{22}}\right) \sigma\left(\frac{x_{2}}{\gamma_{22}}\right)$	$\left(\frac{x_2 - b_{22}}{\gamma_{22}}\right)$ $\left(\frac{-b_{21}}{\gamma_{21}}\right) \sigma\left(\frac{x_2 - c_{22}}{\gamma_{22}}\right)$	$\left(\frac{b_{22}}{2}\right)$		
	$+ \beta_{22} \left(\sigma \left(\frac{x_1 - b_{11}}{\gamma_{11}} \right) - \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \right) \left(\sigma \left(\frac{x_2 - b_{12}}{\gamma_{12}} \right) - \sigma \left(\frac{x_2 - b_{22}}{\gamma_{22}} \right) \right) \\ + \beta_{32} \sigma \left(\frac{x_1 - b_{21}}{\gamma_{22}} \right) \left(\sigma \left(\frac{x_2 - b_{12}}{\gamma_{22}} \right) - \sigma \left(\frac{x_2 - b_{22}}{\gamma_{22}} \right) \right)$								
		+ /	$\beta_{12}\left(1-a\right)$	$\sigma\left(\frac{x_1-b}{\gamma_{11}}\right)$	$\left(\frac{b_{11}}{c}\right)$	$\sigma\left(\frac{x_2 - b_{12}}{\gamma_{12}}\right) - \sigma\left(\frac{x_2 - b_{12}}{\gamma_{12}}\right) $	$\left(\frac{x_2 - \gamma_2}{\gamma_2}\right)$	$\left(\frac{b_{22}}{2}\right)$	
		+ /	$\beta_{31}\sigma\left(\frac{x_1}{x_1}\right)$	$\left(\frac{\gamma_{11}}{\gamma_{21}}\right)$	$\left(1 - \sigma \right) \left(\frac{x}{z}\right)$	$\frac{\gamma_{21}}{\gamma_{22}} \right) \left(\left(\frac{1}{\gamma_{12}} \right) \right)$	γ1:	2 //	
	$f_{\mathcal{E}}(x)$	$(x_1, x_2) = (x_1, x_2) + (x_2, x_2)$	$\beta_{11} \left(1 - \sigma \right)$ $\beta_{21} \left(\sigma \left(\frac{x}{-\tau} \right) \right)$	$\sigma\left(\frac{x_1-\gamma_{11}}{\gamma_{11}}\right)$	$\left(\frac{b_{11}}{a}\right)\left(1\right)$	$\frac{1-\sigma\left(\frac{x_2-b_{12}}{\gamma_{12}}\right)}{\frac{1-b_{21}}{\gamma_{12}}}\right)\left(1-\sigma\left(\frac{1-\sigma_{12}}{\gamma_{12}}\right)\right)$	$(x_2 - x_2 - x_2)$	$\left(\frac{b_{12}}{b_{12}}\right)$	
	r ()		0 (1	$(x_1 - $	$b_{11} \rangle \rangle \langle 1 \rangle$	$(x_2 - b_{12})$			

Figure 5: Decomposition $f_{\mathcal{E}}$ into f_1 and f_2 . Each cell represents the support created by $\{b_{12}, b_{22}, b_{12}, b_{21}\}$ and β_{ij} is the height parameter corresponding to each cell.

 f_1

 ϕ_3

 ϕ_2

 ϕ_1

 f_2

Step 1) Decomposition of $f_{\mathcal{E}}$. As shown in Figure 5, $f_{\mathcal{E}}$ can be decomposed into f_1 and f_2 which satisfy the sum-to-zero condition, and f_1 and f_2 are defined as

$$f_{1}(x_{1}, x_{2}) = (\beta_{11} - \alpha_{1}) \left(1 - \sigma \left(\frac{x_{1} - b_{11}}{\gamma_{11}} \right) \right) \left(1 - \sigma \left(\frac{x_{2} - b_{12}}{\gamma_{12}} \right) \right) + \eta_{1} \sigma \left(\frac{x_{1} - b_{11}}{\gamma_{11}} \right) \left(1 - \sigma \left(\frac{x_{2} - b_{12}}{\gamma_{12}} \right) \right)$$

- $+ (\beta_{12} \alpha_2) \left(1 \sigma \left(\frac{x_1 b_{11}}{\gamma_{11}} \right) \right) \left(\sigma \left(\frac{x_2 b_{12}}{\gamma_{12}} \right) \sigma \left(\frac{x_2 b_{22}}{\gamma_{22}} \right) \right)$ $+ \eta_2 \sigma \left(\frac{x_1 b_{11}}{\gamma_{11}} \right) \left(\sigma \left(\frac{x_2 b_{12}}{\gamma_{12}} \right) \sigma \left(\frac{x_2 b_{22}}{\gamma_{22}} \right) \right)$ $+ (\beta_{13} \alpha_3) \left(1 \sigma \left(\frac{x_1 b_{11}}{\gamma_{11}} \right) \right) \sigma \left(\frac{x_2 b_{22}}{\gamma_{22}} \right)$

1078
$$+ (\beta_{13} - \alpha_3) \left(1 - \sigma \left(\frac{\gamma_{11}}{\gamma_{11}} \right) \right) \sigma \left(- \frac{\gamma_{12}}{\gamma_{11}} \right) \sigma \left(\frac{\gamma_{12}}{\gamma_{11}} \right) \sigma \left(\frac{\gamma_{12}}{\gamma_{11}} \right) \sigma \left(\frac{\gamma_{12}}{\gamma_{12}} \right) \sigma \left(\frac{\gamma$$

 $+\eta_3\sigma\bigg(\frac{x_1-b_{11}}{\gamma_{11}}\bigg)\sigma\bigg(\frac{x_2-b_{22}}{\gamma_{22}}\bigg)$

 $f_{\mathcal{E}}$



Figure 6: **Decomposition** f_2 into f_{21} and f_{22} . Each cell represents the support created by $\{b_{12}, b_{22}, b_{21}\}$ and β_{ij} is the height parameter corresponding to each cell.

Step 2) Decomposition of f_2 . Similar to step 1, f_2 can be decomposed into f_{22} and f_{12} , as described in Figure 6. Note that f_{22} and f_{12} are ANOVA-NODTs which are defined as

$$\begin{aligned} & 1112 \\ & f_{21}(x_1, x_2) = \tau_{11} \left(1 - \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \right) \left(1 - \sigma \left(\frac{x_2 - b_{22}}{\gamma_{22}} \right) \right) \\ & + \tau_{21} \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \left(1 - \sigma \left(\frac{x_2 - b_{22}}{\gamma_{22}} \right) \right) \\ & 1115 \\ & + \tau_{12} \left(1 - \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \right) \sigma \left(\frac{x_2 - b_{22}}{\gamma_{22}} \right) \end{aligned}$$

1116
$$(\gamma_{21}) (\gamma_{2})$$

1117 $(x_1 - b_{21}) (x_2 - b_{22})$

1109

1132

1118 $+ \tau_{22}\sigma\left(\frac{x_1 - b_{21}}{\gamma_{21}}\right)\sigma\left(\frac{x_2 - b_{22}}{\gamma_{22}}\right)$

1119 1120 where $\tau_{11} = -\mathbb{E}\left[\sigma\left(\frac{X_2-b_{22}}{\gamma_{22}}\right)\right](\alpha_3-\alpha_2)$ and $\tau_{12}, \tau_{21}, \tau_{22}$ are uniquely determined by sum-to-zero 1121 condition, and

$$f_{22}(x_1, x_2) = h_{11} \left(1 - \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \right) \left(1 - \sigma \left(\frac{x_2 - b_{12}}{\gamma_{12}} \right) \right) + h_{21} \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \left(1 - \sigma \left(\frac{x_2 - b_{12}}{\gamma_{12}} \right) \right) + h_{12} \left(1 - \sigma \left(\frac{x_1 - b_{21}}{\gamma_{21}} \right) \right) \sigma \left(\frac{x_2 - b_{12}}{\gamma_{12}} \right)$$

1127
1128
$$+ h_{22}\sigma\left(\frac{x_1 - b_{21}}{\gamma_{21}}\right)\sigma\left(\frac{x_2 - b_{12}}{\gamma_{12}}\right)$$

1129 where $h_{11} = \alpha_1 - \tau_{11}$ and h_{12}, h_{21}, h_{22} are uniquely determined by sum-to-zero condition. Also, 1130 another function f_1 can be decomposed into two ANOVA-NODTs in the same way as f_2 is decom-1131 posed.

1133 Case of $d \ge 3$. Similar to the toy example at d = 2, the ensemble function $f_{\mathcal{E}}$ can also be decomposed into ANOVA-NODTs when $d \ge 3$.

DETAILS FOR ANOVA-NODT В

In this section, we describes the specific form of ANOVA-NODT h^S for $S \subseteq [p]$.

B.1 ANOVA-NODT FOR THE MAIN EFFECT.

Without loss of generality, we assume that $S = \{i\}$. We can express ANOVA-NODT as:

$$h^{i}(x_{i}) = \beta_{1,i}c_{1}(x_{i}|b_{1,i},\gamma_{1,i}) + \beta_{2,i}c_{2}(x_{i}|b_{1,i},\gamma_{1,i}),$$

$$(7)$$

where $(\beta_{1,i}, \beta_{2,i})'$ are the height parameters of the terminal nodes,

$$c_1(x|b_{1,i},\gamma_{1,i}) = 1 - entmax_{\nu} \left(\left(\frac{x - b_{1,i}}{\gamma_{1,i}}, 0 \right)' \right)_1.$$

$$c_2(x|b_{1,i},\gamma_{1,i}) = 1 - c_1(x|b_{1,i},\gamma_{1,i})$$

and $b_{1,i}$ and $\gamma_{1,i}$ are learnable parameters. To satisfy the sum-to-zero condition, we require

$$\mathbb{E}[h^{i}(X_{i})] = \beta_{1,i}\mathbb{E}[c_{1}(X_{i}|b_{1,i},\gamma_{1,i})] + \beta_{2,i}\mathbb{E}[c_{2}(X_{i}|b_{1,i},\gamma_{1,i})] = 0$$

Without loss of generality, for a given $\theta_i \in \mathbb{R}$, let $\beta_{1,i} = \theta_i$. Then, we have $\beta_{2,i} = -I_{1,i}^{(b_{1,i},\gamma_{1,i})}\theta_i$, where $I_{1,i}^{(b_{1,i},\gamma_{1,i})} = -\frac{\mathbb{E}[c_1(X_i|b_{1,i},\gamma_{1,i})]}{\mathbb{E}[c_2(X_i|b_{1,i},\gamma_{1,i})]}$

Therefore, the equation (7) is expressed as

$$h^{i}(x_{i}) = \theta_{i}c_{1}(x_{i}|b_{1,i},\gamma_{1,i}) + \theta_{i}I_{1,i}^{(b_{1,i},\gamma_{1,i})}c_{2}(x_{i}|b_{1,i},\gamma_{1,i}),$$

where $\theta_i, b_{1,i}$ and $\gamma_{1,i}$ are learnable free parameters.

Monotone constraint. $c_1(x_i|b_{1,i},\gamma_{1,i})$ is a decreasing function with respect to x_i , and accord-ingly, $c_2(x_i|b_{1,i}, \gamma_{1,i})$ is an increasing function with respect to x_i . Therefore, if θ_i is greater than 0, then $h^i(x_i)$ becomes a decreasing function with respect to x_i since $I_{1,i}^{(b_{1,i},\gamma_{1,i})}$ is always less than 0.

B.2 ANOVA-NODT FOR SECOND ORDER INTERACTION.

Without loss of generality, we consider the component $S = \{i, k\}$. For simplicity, we denote $c_h(x_i|b_{z,(i,k)}, \gamma_{z,(i,k)})$ as $c_{h,z}(x_i)$. Then, for the component S, ANOVA-NODT can be expressed as: (. 1)

$$h^{(i,k)}(x_i, x_k) = \beta_{(1,1)',(i,k)} c_{1,1}(x_i) c_{1,2}(x_k)$$

$$+\beta_{(1,2)',(i,k)}c_{1,1}(x_i)c_{2,2}(x_k)$$

1172
$$+ \beta_{(2,1)',(i,k)}c_{2,1}(x_i)c_{1,2}(x_k)$$

 $+ \beta_{(2,2)',(i,k)}c_{2,1}(x_i)c_{2,2}(x_k),$

where $(\beta_{(1,1)',(i,k)},\beta_{(1,2)',(i,k)},\beta_{(2,1)',(i,k)},\beta_{(2,2)',(i,k)})$ are the height parameter vector of the ter-minal nodes and

1177
1178
1179
1180

$$c_{1,t}(x) = 1 - entmax_{\nu} \left(\left(\frac{x - b_{t,(i,k)}}{\gamma_{t,(i,k)}}, 0 \right)' \right)_{1}$$

for
$$t = 1, 2$$
. To satisfy sum-to-zero condition, $h^{i,k}$ has to meet the followings: for $x_k \in \mathcal{X}_k$,

1182
$$\mathbb{E}[h^{(i,k)}(X_i, x_k)] = \beta_{(1,1)',(i,k)}c_{1,2}(x_k)\mathbb{E}[c_{1,1}(X_i)]$$

$$+ \beta_{(1,2)',(i,k)} c_{2,2}(x_k) \mathbb{E}[c_{1,1}(X_i)]$$

1185
$$+ \beta_{(2,1)',(i,k)} c_{1,2}(x_k) \mathbb{E}[c_{2,1}(X_i)]$$

1186
$$+ \beta_{(2,2)',(i,k)} c_{2,2}(x_k) \mathbb{E}[c_{2,1}(X_i)]$$

1187 $- \beta_{(2,2)',(i,k)} c_{2,2}(x_k) \mathbb{E}[c_{2,1}(X_i)]$

$$= 0$$

1188	and for $x_i \in \mathcal{X}_i$,
1189	$\mathbb{E}[h^{(i,k)}(x, Y_{i})] = \beta_{(i,k)}(x, y_{i}) \mathbb{E}[a, a(Y_{i})]$
1101	$\mathbb{E}[n^{k-1}(x_{i}, X_{k})] = \beta(1, 1)^{j}, (i, k) \in [1, 1(X_{i}) \mathbb{E}[0, 2(X_{k})]]$
1102	$+\beta_{(1,2)',(i,k)}c_{1,1}(x_i)\mathbb{E}[c_{2,2}(X_k)]$
1193	$+ \beta_{(2,1)',(i,k)} c_{2,1}(x_i) \mathbb{E}[c_{1,2}(X_k)]$
1194	$+ \beta_{(2,2)',(i,k)}c_{2,1}(x_i)\mathbb{E}[c_{2,2}(X_k)]$
1195	= 0.
1196	The following conditions are a rephrase of the above conditions:
1197	$\beta_{(2,1)}(x) \mathbb{E}[c_{1,1}(X)] + \beta_{(2,1)}(x) \mathbb{E}[c_{2,1}(X)] = 0$
1190	$\mathcal{P}_{(1,1)',(i,k)} \boxtimes [c_{1,1}(\mathcal{X}_{i})] + \mathcal{P}_{(2,1)',(i,k)} \boxtimes [c_{2,1}(\mathcal{X}_{i})] = 0$
1200	$\beta_{(1,2)',(i,k)}\mathbb{E}[c_{1,1}(X_i)] + \beta_{(2,2)',(i,k)}\mathbb{E}[c_{2,1}(X_i)] = 0$
1201	$\beta_{(1,1)',(i,k)}\mathbb{E}[c_{1,2}(X_k)] + \beta_{(1,2)',(i,k)}\mathbb{E}[c_{2,2}(X_k)] = 0$
1202	$\beta_{(2,1)',(i,k)}\mathbb{E}[c_{1,2}(X_k)] + \beta_{(2,2)',(i,k)}\mathbb{E}[c_{2,2}(X_k)] = 0$
1203 1204	For $\theta_{i,k} \in \mathbb{R}$, let $\beta_{(1,1)',(i,k)} = \theta_{i,k}$. Then, ANOVA-NODT $h^{(i,k)}(x_i, x_k)$ can be expressed as
1205 1206	$h^{(i,k)}(x_i, x_k) = \theta_{i,k} c_{1,1}(x_i) c_{1,2}(x_k)$
1207	$+ \theta_{i,k} I_{1,(i,k)} c_{1,1}(x_i) c_{2,2}(x_k)$
1208	$+ \theta_{i k} I_{2(i k)} c_{21}(x_i) c_{12}(x_k)$
1209	$+ \theta_{i,k} I_{2,(i,k)} 2, 1(x_i) + 1, 2(x_k)$
1210	$+ v_{i,k} r_{3,(i,k)} c_{2,1}(x_i) c_{2,2}(x_k)$
1211	where
1212	$\mathbb{E}[c_{1,2}(X_k)]$
1213	$T_{1,(i,k)} = -rac{\mathbb{E}[c_{2,2}(X_k)]}{\mathbb{E}[c_{2,2}(X_k)]}$
1214	$\mathbb{E}[c_{1,1}(X_i)]$
1215	$I_{2,(i,k)} = -\frac{\mathbb{E}[c_{2,1}(X_i)]}{\mathbb{E}[c_{2,1}(X_i)]}$
1210	$I_{3(i,k)} = I_{1(i,k)} I_{2(i,k)}.$
1218	-3,(i,k) $-1,(i,k)$ $-2,(i,k)$
1219	and $\theta_{i,k}$ as well as $b_{1,(i,k)}, \gamma_{1,(i,k)}$ and $b_{2,(i,k)}, \gamma_{2,(i,k)}$ are learnable free parameters.
1220	
1221	
1222	
1223	
1224	
1225	
1220	
1227	
1229	
1230	
1231	
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1241	

¹²⁴² C DETAILS OF THE EXPERIMENTS

All models except XGB were trained via Adam optimizer. Likewise in Popov et al. (2019), we set $\nu = 1.5$ for $entmax_{\nu}$ in all the experiments. All experiments are run with RTX 3090, RTX 4090, and 24GB memory.

C.1 DETAILS FOR SYNTHETIC DATASETS

Table 6: Test suite of synthetic functions.

$f^{(1)}$	$Y = 10X_1 + 10X_2 + 20(X_3 - 0.3)(X_3 - 0.6) + 20X_4 + 5X_5 + 10\sin(\pi X_1 X_2) + \epsilon$
$f^{(2)}$	$Y = \pi^{X_1 X_2} \sqrt{2X_3} - \sin^{-1}(X_4) + \log(X_3 + X_5) - \frac{X_9}{X_{10}} \sqrt{\frac{X_7}{X_8}} - X_2 X_7 + \epsilon$
$f^{(3)}$	$Y = \exp X_1 - X_2 + X_2X_3 - X_3^{2 X_4 } + \log(X_4^2 + X_5^2 + X_7^2 + X_8^2) + X_9 + \frac{1}{1 + X_{10}^2} + \epsilon$

Table 7: Distribution of input features in synthetic functions.

$f^{(1)}$	$X_1, X_2, X_3, X_4, X_5 \sim^{iid} U(0, 1)$
$f^{(2)}$	$X_1, X_2, X_3, X_6, X_7, X_9 \sim^{iid} U(0, 1) \text{ and } X_4, X_5, X_8, X_{10} \sim^{iid} U(0.6, 1).$
$f^{(3)}$	$X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10} \sim^{iid} U(-1, 1)$

The synthetic function $f^{(1)}$ is a slightly modified version of Friedman's synthetic function used in Chipman et al. (2010). $f^{(2)}$ and $f^{(3)}$ are taken from the synthetic functions used in the interaction detection experiments in Tsang et al. (2017). We generate 15K data samples from the distribution in the Table 7 and functions in the Table 6. Also, we divide them into train, validation and test datasets with ratio 0.7, 0.1 and 0.2, respectively. For all of the synthetic functions, the number of trees for component *S*, K_S , is set to 30.

C.2 DETAILS OF THE EXPERIMENTS WITH REAL DATASETS.

Table 8:	Descriptions	of real	datasets.
----------	--------------	---------	-----------

		-		
Dataset	Size	Number of features	Problem	Number of Class
CALHOUSING	21k	8	Regression	-
WINE	4k	11	Regression	-
ONLINE	40k	58	Regression	-
ABALONE	4k	10	Regression	-
FICO	10k	23	Classification	2
CHURN	7k	39	Classification	2
Credit	284k	30	Classification	2
Letter	20k	16	Classification	2
Drybean	13k	16	Classification	7
MICROSOFT	960k	136	Regression	-
Үаноо	700k	699	Regression	-
MADELON	2.6k	500	Classification	2
CelebA	200K		Classification	2

Dataset	ANOVA-N ¹ ODE	ANOVA-N ² ODE
CALHOUSING	10	10
WINE	100	10
Online	10	10
Abalone	50	10
FICO	30	30
CHURN	10	10
Credit	10	5
Letter	50	10
Drybean	50	100
MICROSOFT	10	10
Үаноо	10	10
MADELON	10	40

Table 9: K_S in ANOVA-N¹ODE and ANOVA-N²ODE

1309

1296

Implementation of baseline model. We conduct experiments for all baseline models (NAM, NBM, NODE-GAM, NODE) using the official source code. For XGB, we utilize the xgboost package (Chen & Guestrin, 2016) for our experiments.

1314
 1315
 1316
 Data descriptions. Table 8 summarizes the descriptions of 9 real datasets we analyze in the numerical studies.

Data preprocessing. Minimax scaling is applied to NAM and NBM, while standard scaling was used for NODE-GAM, NODE, and XGB. For ANOVA-NODE, transformation using quantiles of a uniform distribution is performed to satisfy sum-to-zero condition stably during training. Addition ally, all categorical features is encoded using one-hot encoding.

1321

Learning rate. For all models except XGB, we set the learning rate of Adam optimizer as 5e-3
and batch size is 4096. We find the optimal learning rate of XGB via grid search.

1324

Model hyperparameters. Table 9 presents the number of NODTs K_S used in ANOVA-N¹ODE and ANOVA-N²ODE for real datasets. In NA¹M, the dimensions of the hidden layers of each component consists of [64,32,16] for MICROSOFT, YAHOO and MADELON, and [64,64,32] for other datasets. In NA²M, the hidden layers consist of [64,32,16] for the ONLINE, CREDIT and DRYBEAN datasets, [64,16,8] for MICROSOFT, YAHOO and MADELON, and [64,64,32] for the other datasets.

For XGB, NODE, and NODE-GAM, we randomly split the train, validation and test data into the ratio 70/10/20 and evaluated its performance on the validation dataset using the model trained on the train dataset. We repeated this process 10 times with randomly split data, resulting in 10 prediction performance values for the validation datasets. Then, we selected the optimal hyper-parameters by the grid search based on the average of the prediction performance values for the validation datasets.

Finally, with the optimal hyper-parameters selected earlier, we fixed the model's hyper-parameters and used the 10 train-test dataset pairs obtained from the previous data splitting to train the model on the train datasets and evaluate its performance on the test datasets.

- For XGB, the range of hyper-parameters for the grid search is as below.
- 1340

1347

- The number of tree : {50,100,200,300,400,500,600,700,800,900,1000}
- max depth : {3, 5, 7}
 - learning rate : {0.0001, 0.005, 0.01, 0.05, 0.1}

For NODE, NODE-GA¹M and NODE-GA²M, the range of hyper-parameters for the grid search is as below.

- The number of layer : {2, 4, 8}
- tree depth : {6, 8}
 - The number of trees in each layer : {256, 512}

Selected components by NID for high-dimensional datasets. Table 10 presents the number of components used in training ANOVA-N²ODE and baseline models. All main effects are used, and the second order interactions are selected using NID (Tsang et al., 2017). For MICROSOFT, 300 second order interactions are used; 500 for YAHOO, 500; and 300 for MADELON.

Table 10: The number of components used in training ANOVA-N²ODE, NA²M, and NB²M.

	•	0	· · ·
Dataset	MICROSOFT	Үаноо	MADELON
# of selected components	136(Main) + 300(2nd)	699(Main) + 500(2nd)	500(Main) + 300(2nd)

1360 C.3 EXPERIMENT DETAILS FOR IMAGE DATASET.

For CELEBA image dataset, we use the joint bottleneck model in Koh et al. (2020). The main idea of the joint bottleneck model (Koh et al., 2020) is not to directly input the embedding vector derived from image data through a CNN into a classifier for classification. Instead, CNN predicts given concepts (attributes) for the image, and the predicted values for these concepts are then used as an input of classifier. In Koh et al. (2020), they used DNN as a classifier which is a black box model. In this paper, we replace DNN with ANOVA-NODE, NAM, NBM and NODE-GAM. For CNN, we linear heads on the bottom of the pretrained ResNet18.

All models are trained via the Adam optimizer with a 1e-3 learning rate and the batch size for training is 256. In ANOVA-N¹ODE and ANOVA-N²ODE, K_S is set to 10 and 3, respectively. In NA¹M and NA²M, we use a neural network consisting of 3 hidden layer with sizes (64, 16,8) and (64,4,2), respectively. In NB¹M and NB²M, we use 100 basis neural networks consisting of 3 hidden layer with sizes (256,128,128) and (128,64,64) for the basis model, respectively. Also, for NODE-GA¹M and NODE-GA²M, the number of trees is set to 125 and 50, respectively, and the depth and the number of layers are set to be 6 and 4, respectively.

The number of trees for each component

ANOVA-N¹ODE

ANOVA-N²ODE

1404 D ABLATION STUDIES. 1405

1406

1407 1408

1409 1410 1411

1420

1421

1422

1423

1430

1439 1440 D.1 THE CHOICE OF THE NUMBER OF TREES IN ANOVA-NODE.

Table 11 presents the averages and standard deviations of prediction performance of ANOVA-NODE 1412 based on 10 randomly sampled data of ABALONE datasets for various values of the number of trees 1413 K_S . We observe that K_S around 50 yields the best results for ANOVA-N¹ODE and K_S around 10 1414 for ANOVA- N^2 ODE. We report that similar results are obtained for other datasets. 1415

Table 11: Results of prediction performance for various numbers of trees.

2.163 (0.08)

2.102 (0.08)

10

2.160 (0.09)

2.087 (0.08)

50

2.135 (0.09)

2.105 (0.08)

100

2.159 (0.08)

2.122 (0.08)

1424		
1425		
1426		
1427		
1428	D 2	IMPACT OF THE INITIAL VALUES OF MODEL PARAMETERS TO STABILITY
1429	D.2	

2.176 (0.09)

2.103 (0.08)

1431 We investigate how the choice of initial values of the model parameters affects the stability of the es-1432 timated components by ANOVA-NODE, NAM and NBM by analyzing a synthetic dataset generated 1433 from $f^{(1)}$. 1434

We conducted 5 trials on the same train/test/validation dataset, and the results are presented in Figure 1435 7 ,8 and 9. We observe that NA²M and NB²M frequently estimate the true function inaccurately. 1436 In contrast, ANOVA-N²ODE consistently estimates the components accurately regardless of the 1437 choice of the initial values. 1438



Figure 7: Plots of the functional relations of the main effects in ANOVA-N²ODE on synthetic datasets generated from $f^{(1)}$.



¹⁵¹² E ILLUSTRATION OF INTERPRETABILITY OF ANOVA-NODE.

1514 E.1 ILLUSTRATION OF INTERPRETABILITY

1516 1517

We consider the two concepts of interpretation: Local and Global which are roughly defined as:

Local Interpretation: Information about how each feature of a given datum affects the prediction.
SHAP is a notable example of local interpretation. For the functional ANOVA model, the predictive values of each component at a given datum would be considered as local interpretation.

Global Interpretation: Information about how each feature is related to the final prediction model. The importance scores of each feature (e.g. global SHAP (Molnar, 2020)) and the functional relations between each feature and the prediction model (e.g. the dependency plot of SHAP (Molnar, 2020) are examples of global interpretation. For the functional ANOVA model, the importance score, which can be defined as the l_1 norm of the corresponding component as is done in Section 4.2, and the functional relation identified by the functional form of each component are two tools for global interpretation.

1529 E.1.1 ILLUSTRATION OF INTERPRETABILITY ON CALHOUSING DATASET.

 Table 12: Feature descriptions of CALHOUSING dataset.

Feature name	Index	Description	Feature type
MedInc	1	Median income in block	Numerical
HouseAge	2	Median house age in block	Numerical
AveRooms	3	Average number of rooms	Numerical
AveBedrms	4	Average number of bedrooms	Numerical
Population	5	Population in block	Numerical
AveOccup	6	Average house occupancy	Numerical
Latitude	7	Latitude of house block	Numerical
Longitude	8	Longitude of house block	Numerical

Local Interpretation on CALHOUSING dataset. We conduct an experiment on CALHOUSING (Pedregosa et al., 2011a) dataset to illustrate local interpretation of ANOVA-N¹ODE. Note that ANOVA-N¹ODE is given as

1542 1543

1541

1531

1544

1546

1550

1551

1553

1554

1557

 $\hat{f}_{\text{ANOVA-N}^1\text{ODE}}(\mathbf{x}) = \sum_{j=1}^8 \hat{f}_j(x_j).$

Thus, it is reasonable to treat $\hat{f}_j(x_j)$ as the contribution of x_j to $\hat{f}(\mathbf{x})$. In fact, we have seen in Section 3.1 that this contribution is equal to SHAP (Lundberg, 2017). As an illustration, for a given datum

 $\mathbf{x} = (-0.2378, -0.4450, 0.0036, -0.1531, 0.3814, -0.067, 0.5541, -0.1111)^{\top},$

the contributions of each feature to $\hat{f}(\mathbf{x})$ are

 $(\hat{f}_1, ..., \hat{f}_8) = (-4.9900, 0.3278, -0.0456, 0.4432, -0.1730, 2.7521, -11.6190, 6.5184).$

That is, the 7th variable contributes most to the prediction value of $\hat{f}(\mathbf{x})$, which can be interpreted as 'the housing price is low because the latitude is not good'.

Global Interpretation on CALHOUSING dataset. Figure 10 and Table 13 present the functional relations of each input feature to the prediction model learned by ANOVA-NO¹DE and their importance scores. From these results, we can see that the location is the most important features and the housing price on the south-west area is the most expensive.

1562Table 14 describes the 10 most important components with descending order of the importance1563scores of ANOVA-NO²DE normalized by the maximum importance score. The results are bit differ-1564ent from those of ANOVA-NO¹DE. In particular, the interaction between 'latitude' and 'longitude'1565emerges as a new important feature while the main effects of 'latitude' and 'longitude' become lessimportant.



1623

1624 1625

1626

1673

1627	serve th	at the predic	ction perfo	rmances o	f ANOVA-Ì	N^1 ODE and Al	NOVA-N ² OI	DE are con	nparable or
1628	superior	r to their co	mpetitors.						
1629									
1630	Attribu	tes to whic	h monoto	ne constra	ints are ap	plied. For at	tributes 'Bal	d', 'Big No	ose', 'Goa-
1631	tee' and	'Mustache	', we apply	the increa	asing monot	one constraint	, while for at	tributes 'A	rched Eye-
1632	brows',	'Attractive'	, 'Heavy M	Makeup', '	No Beard',	'Wearing Earn	rings', 'Wear	ing Lipstic	ck', 'Wear-
1633	ing Nec	klace', 'We	aring Necl	ctie', we u	sed the deci	easing monoto	one constrain	ıt.	
1634			Table 16.	Importon	a coore for	the 2 importan	taamnanant	-	
1635				Importanc	e scores for	the 5 Importan	t components	š.	
1636			Components	Monotone	No Beard	Wearing Lipstick	Heavy Makeup	<u>,</u>	
1637			Score	0	0.794	0.465	0.210	_	
1638									
1639									
1640		Table 17:	Results of	local inter	pretation wit	h and without	the monotone	e constrain	t.
1641		Image index	Monotone	Heavy Make	eup No beard	Wearing Lipstic	k classified lat	oel True lab	oel
1642		2-1	X	0.030	0.035	0.093	male	female	2
16/12		2-1	0	-0.080	-0.161	-0.106	female	female	<u> </u>
1643		2-2	<u>л</u> О	-0.081	-0.183	-0.106	female	male	
1044			_						
1045									
1040				1000					
1647				Cartha M		105			
1648				P. May Ser	1	CHERN WAY	AN CT		
1649					1		1253		
1650			13	1000	11		A DA THE		
1651			1	1 U	18 The second	Aller	MINS _		
1652			10		180	Ser and			
1653			91			100			
1654							1.000		
1655				Image 2-7		Image 2-2			
1656				Figure 1	12: Misclass	ified two image	es.		
1657									
1658									
1659		,	Table 18: A	Accuracies	(standard d	eviations) on C	ELEBA datas	et.	
1660	ANOVA-	$-N^{1}ODE NOE$	$\frac{\text{DE-GA}^{1}\text{M}}{1(0.006)} = 0$	NA ¹ M	NB ¹ M	ANOVA-N ² ODE	NODE-GA ² M	NA ² M	NB ² M
1661	0.985 (0.001) 0.98	1 (0.000) 0.	982 (0.002)	0.980 (0.002)	0.980 (0.001)	0.981 (0.000)	0.960 (0.001)	0.980 (0.002)
1662									
1663									
1664									
1665									
1666									
1667									
1668									
1660									
1670									
10/0									
16/1									
1672									

1620	Table 15: Result	s of prediction	performa	nce of A	NOVA-NODE	with an	d without th	ne monotone	con-
1621	straint.								
1622				Measure	ANOVA-N ¹	DDE AN	OVA-N ² ODE		

Accuracy ↑

Comparison with baseline models in terms of prediction performance. In Table 18, we ob-

Without Monotone constraint Accuracy ↑

With Monotone constraint

0.985 (0.001)

0.984 (0.001)

ANOVA-N²ODE

0.986 (0.001)

0.985 (0.001)

F ADDITIONAL EXPERIMENTS FOR THE STABILITY OF ANOVA-NODE.

F.1 STABILITY OF THE ESTIMATED COMPONENTS ON VARIATIONS OF TRAINING DATA

We investigate the stability of components estimated by ANOVA- N^1 ODE and ANOVA- N^2 ODE when training data vary. We use CALHOUSING (Pedregosa et al., 2011a) and WINE (Cortez et al., 2009) datasets and compare ANOVA-NODE with NAM and NBM.

Experiment for CALHOUSING dataset. Figures 13, 14 and 15 present the plots of the functional relations of the main effects estimated by ANOVA-N¹ODE, NA¹M, and NB¹M for 5 randomly sam-pled training datasets. Figures 13, 14 and 15 present the plots of the functional relations of the main effects estimated by ANOVA-N²ODE, NA²M, and NB²M for 5 randomly sampled training datasets. We observe that the 5 main components estimated by ANOVA-N¹ODE and ANOVA-N²ODE are relatively much more stable compared to NAM and NBM. Note that as seen in Figure 17, we observe that in NA²M, some components are estimated as a constant function, which would be partly because the main effects are absorbed into the second order interactions.

Experiment for WINE dataset. Figures 19, 20 and 21 present the plots of the functional relations of the main effects estimated by ANOVA-N¹ODE, NA¹M, and NB¹M for 5 randomly sampled training datasets. Figures 22, 23 and 24 present the plots of the functional relations of the main effects estimated by ANOVA-N²ODE, NA²M, and NB²M for 5 randomly sampled training datasets. The results are similar to those of CALHOUSING dataset.



Figure 13: Plots of the functional relations of the main effects in ANOVA-N¹ODE on CALHOUSING dataset.







Figure 22: Plots of the functional relations of the main effects in ANOVA-N²ODE on WINE dataset.



F.2 STABILITY OF ANOVA-SHAP ON CALHOUSING DATASET

1926We conduct an experiment to evaluate the stability of ANOVA-SHAP on CALHOUSING dataset.1927We calculate the global importance of features for 10 trials using the l_1 norm of the ANOVA-1928SHAP values defined in (2). Finally, we compute the stability score of ANOVA-SHAP as the av-1929erage of Hamming distance between the global importance ranks across all pairs in the trials. Table193019 presents the results of stability scores of ANOVA-SHAP which are normalized by the that of1931ANOVA-N¹ODE (ANOVA-N²ODE). We confirm that our model provides significantly more stable1932ANOVA-SHAP interpretations compared to other baseline models.

1933

1934 Table 19: Results of average of Hamming distance. A smaller distance indicates that the interpretation of ANOVA-SHAP is more stable.

Model	ANOVA-N ¹ ODE	NA ¹ M	NB ¹ M
Average of Hamming distance	1.000	6.188	2.408
Model	ANOVA-N ² ODE	NA ² M	NB ² N
Average of Hamming distance	1.000	5.157	2.663

1938 1939

1936 1937

1941 F.3 SCATTER PLOTS OF STABILITY SCORE

1942 In this section, we present the scatter plots of the stability score on WINE dataset. It is obvious that 1943 ANOVA N¹ODE as well as ANOVA N²ODE are more stable in estimation of the components than NAM and NBM.



Figure 25: Scatter plots of the stability scores on WINE dataset. Figures (a) and (d) are the scatter plots for the stability scores of the main effects, where the x-axis is the stability score of ANOVA-N¹ODE and the y-axis is the stability scores of NA¹M in (a) and NB¹M in (d). Figures (b) and (e) are the scatter plots for the stability scores of the main effects, where the x-axis is the stability score of ANOVA-N²ODE, and the y-axis is the stability score of NA^2M in (b) and NB^2M in (e). Figures (c) and (f) compare the stability scores of the second order interactions of ANOVA-N²ODE to those of NA²M and NB²M. Each dot in the scatter plots corresponds to each component.

G ADDITIONAL EXPERIMENTS ON HIGH-DIMENSIONAL DATASETS

G.1 RESULTS OF UNNORMALIZED STABILITY SCORE.

Table 20 presents the original stability scores SC(f)/|S| of the normalized stability scores presented in Table 1.

2004 2005 2006

2003

1998

2000

2001 2002

2007 2008 2009

2022 2023 2024

Table 20: Results of stability scores in each model on the real datasets.

Dataset	ANOVA-N ¹ ODE	NA ¹ M	NB ¹ M	ANOVA-N ² ODE	NA ² M	NB ² M
CALHOUSING	0.012	0.045	0.039	0.035	0.071	0.075
WINE	0.011	0.058	0.043	0.049	0.087	0.065
ONLINE	0.031	0.076	0.054	0.052	0.072	0.072
ABALONE	0.008	0.013	0.026	0.028	0.047	0.038
FICO	0.035	0.046	0.046	0.048	0.089	0.075
CHURN	0.017	0.027	0.047	0.047	0.089	0.080
CREDIT	0.021	0.069	0.025	0.036	0.089	0.053
LETTER	0.017	0.022	0.014	0.026	0.075	0.081
DRYBEAN	0.028	0.074	0.070	0.053	0.088	0.081

2015 2016 G.2 EXTENSION TO NEURAL BASIS MODEL

Similarly to NBM (Radenovic et al., 2022), we can consider extension of ANOVA-NODE using ANOVA-NODTs as basis functions. We call this extension model as NBM-NODE. Consider basis ANOVA-NODTs i.e., $\{h_k(x|\phi_k) : \mathbb{R} \to \mathbb{R}, k = 1, ..., B\}$, then the NBM-N¹ODE $f_{\text{NBM-N}^1\text{ODE}}(\mathbf{x})$ is defined as

$$f_{\text{NBM-N}^1\text{ODE}}(\mathbf{x}) = \sum_{j=1}^p f_{\text{NBM-N}^1\text{ODE}}^j(x_j)w_j \tag{8}$$

where $f_{\text{NBM-N}^1\text{ODE}}^j(x_j) = \sum_{k=1}^B h_k(x_j|\phi_k)a_{jk}$ for j = 1, ..., p and and $h_k(x|\phi_k)$ satisfy the sum-tozero condition with respect to the uniform distribution for μ_j . NBM-N¹ODE can be easily extended to NBM-N²ODE in a similar way as Radenovic et al. (2022).

Figures 26 and 27 show the plots of the functional relations of the main effects estimated by NBM-N¹ODE on the WINE dataset and the CALHOUSING dataset Table 21 shows the prediction performance of NBM-N¹ODE, and Table 22 presents the results of stability scores normalized by the that of ANOVA-N¹ODE. We observe that NBM-N¹ODE also exhibits similar prediction performance and stability to ANOVA-N¹ODE.



Figure 26: Plots of the functional relations of the main effect estimated by NBM-N¹ODE on 5 randomly sampled training data from CALHOUSING dataset.



2106 H ADDITIONAL EXPERIMENTS FOR COMPONENT SELECTION

Table 23 presents the averages and standard deviations of the prediction performance of the models used in the component selection experiment. The three models perform similarly.

2111Table 23: The results of prediction performance. We report the averages and standard deviations of RMSEs2112of ANOVA-N²ODE, NA²M and NB²M on 10 synthetic datasets generated from $f^{(1)}$, $f^{(4)}$ and $f^{(3)}$.2113 $GA^{2}M$

			GA ² M	
Synthetic function	Measure	ANOVA N ² ODE	NA^2M	NB^2M
£(1)	DMSE	3.483	3.474	3.511
J , J	KNISE 4	(0.03)	(0.03)	(0.03)
r(2)	DMCE	0.076	0.088	0.075
J	KNISE ↓	(0.001)	(0.005)	(0.001)
$f^{(3)}$	DMSE	0.161	0.183	0.137
	KNISE ↓	(0.003)	(0.016)	(0.003)

2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2103
2154
2100
2150
21J/

2158 2159

Ι ADDITIONAL EXPERIMENTS FOR PREDICTION PERFORMANCE OF DECISION TREE

Table 24 presents the averages and standard deviations of the prediction performance of decision tree (Breiman, 2017) for 10 trials. We implemented a decision tree by using the scikit-learn python package (Pedregosa et al., 2011b) and turned by using the optuna python package based on below range of hyper-parameters.

- Range of max depth = [2, 12]

- Range of min_samples_leaf = [2,10]
- Range of min_samples_split = [2,10]
- Range of max_leaf_nodes = [2,10]

Table 24: Results of the prediction performance in decision tree and ANOVA-NODE.

2175	Dataset	Measure	Decision Tree	ANOVA-N ¹ ODE	ANOVA-N ² ODE
2176	CALHOUSING	RMSE ↓	0.671 (0.02)	0.614 (0.01)	0.512 (0.01)
2177	WINE	RMSE ↓	0.811 (0.03)	0.725 (0.02)	0.704 (0.02)
2178	ONLINE	RMSE↓	1.119 (0.26)	1.111 (0.25)	1.111 (0.25)
2179	ABALONE	RMSE↓	2.396 (0.08)	2.135 (0.09)	2.087 (0.08)
2190	FICO	AUROC ↑	0.704 (0.02)	0.799 (0.007)	0.800 (0.007)
2100	CHURN	AUROC ↑	0.676 (0.03)	0.839 (0.012)	0.842 (0.012)
2181	Credit	AUROC ↑	0.890 (0.02)	0.983 (0.005)	0.984 (0.006)
2182	LETTER	AUROC †	0.745 (0.001)	0.900 (0.003)	0.984 (0.001)
2183	DRYBEAN	AUROC ↑	0.975 (0.0002)	0.995 (0.001)	0.997 (0.001)
2184					

²²¹⁴ J Additional Experiments for runtime on various datasets.

We conducted additional experiments to assess the improvement in scalability. We consider NA¹M, which has 3 hidden layers with 16, 16, and 8 units; 10 basis DNNs for NB¹M, which has 3 hidden layers with 32, 16, and 16 units; 10 trees for each component in ANOVA-N¹ODE; and 10 basis functions in NBM-N¹ODE. Table 25 presents the results of runtime in NA¹M, NB¹M, ANOVA-N¹ODE, and NBM-N¹ODE on ABALONE, CALHOUSING, and ONLINE datasets. Note that our computational environment consists of RTX 3090 and RTX 4090.

Table 25: Results of runtime in NA¹M, NB¹M, ANOVA-N¹ODE, and NBM-N¹ODE.

Dataset	Size of dataset	# of features	NA ¹ M	NB^1M	ANOVA-N ¹ ODE	NBM-N ¹ ODI
ABALONE	4K	10	6.6 sec	3.0 sec	4.2 sec	1.5 sec
CALHOUSING	21K	8	14.1 sec	4.1 sec	9,7 sec	3.5 sec
ONLINE	40K	58	68 sec	15.6 sec	70 sec	9,8 sec

2268 K ADDITIONAL EXPERIMENTS FOR APPLICABILITY OF ANOVA-SHAP

As explained in Section 3.1, ANOVA-SHAP value of the estimated ANOVA-NODE model can be easily calculated from the estimated components. To evaluate the similarity between ANOVA-SHAP and SHAP without feature independence, we compare ANOVA-SHAP and SHAP values of the estimated ANOVA-N²ODE model on CALHOUSING dataset, where SHAP is computed using Deep-SHAP of Lundberg (2017). Note that the computation time of ANOVA-SHAP was approximately 1,600 times shorter than that of Deep-SHAP.

Figure 28 presents the boxplots of the absolute differences between ANOVA-SHAP and Deep-SHAP values at each data point for the 8 features, based on 1,000 data points which are randomly sampled from the test data, where ANOVA-SHAP values of NAM and NBM are calculated by the equation (2) with the estimated components by NA²M and NB²M, respectively. The absolute differences between ANOVA-SHAP and Deep-SHAP of ANOVA-NODE are distributed around zero which indicates that ANOVA-SHAP is a computationally efficient alternative to Deep-SHAP for ANOVA-NODE. In contrast, the boxplots for NA^2M and NB^2M , which either are far from zero or have large variations in many cases, imply that the formula (2) of ANOVA-SHAP is only applicable when the components satisfy the sum-to-zero condition. The results for stability of ANOVA-SHAP are given in Appendix F.2.



Figure 28: Boxplots of the absolute differences between Deep-SHAP and ANOVA-SHAP values.

²³²² L ANOVA-NODE WITHOUT SUM-TO-ZERO CONDITION

We investigate the performance of ANOVA-NODE without the sum-to-zero condition, which we denote GAM-NODE, by analyzing CALHOUSING and WINE datasets. In GAM-NODE, all of the heights at the terminal nodes of ANOVA-NODE are learnable parameters.

Table 26: **Comparison of ANOVA-NODE and GAM-NODE**. We report the averages of RMSE and stability score (normalized by the that of ANOVA-N¹ODE or NOVA-N²ODE) for 10 trials.

	ANOVA-N ¹ ODE	GAM-N ¹ ODE	ANOVA-N ² ODE	GAM-N ² ODE
CALHOUSING	0.614 (1.000)	0.580 (1.500)	0.512 (1.000)	0.502 (1.690)
WINE	0.725 (1.000)	0.713 (2.550)	0.704 (1.000)	0.690 (1.300)

2333 2334 2335

2324

2328

2330

2331 2332

Table 26 presents (RMSE, stability score) of ANOVA-NODE and GAM-NODE based on 10 randomly selected datasets. Without the sum-to-zero condition, we observe increasing in the stability score. In particular, when the second order interactions are in the model, the main effects are estimated very unstably.

In Table 26, we observe that the prediction performance of GAN-NODE is (slightly) better than that of ANOVA-NODE. One reason could be that ANOVA-NODE is more vulnerable to the local minima problem. Further studies would be worth pursuing.

Figure 29 and 30 present the plots of the functional relations of the main effects on CALHOUSING and WINE dataset in GAM-N²ODE. We observe that GAM-N²ODE estimates the components more unstable compared to ANOVA-N²ODE.

2347 2348

Comparison between NODE-GAM and ANOVA-NODE. In NODE-GAM (Chang et al., 2021), the feature function F^c of NODT at detph c is a sparse weighted sum of input features by using *entmax*_{ν} and temperture parameter T. In other words, for a given depth D and c = 1, ..., D, F^c is defined as below.

$$F^{c}(\mathbf{x}) = \mathbf{x} \cdot entmax_{\nu} \left(\frac{\theta_{F}}{T}\right)$$
$$= \sum_{j=1}^{p} x_{j} w_{j}$$

2358 2359

2360 2361 where $\theta_F = (\theta_{F1}, ..., \theta_{Fp})^{\top}$ is a learnable parameter and $w_j = entmax_{\nu} \left(\frac{\theta_F}{T}\right)_j$. They expect 2362 the projects $[w_j = w_j]$ to be trained as 0 or 1, but there exists the provided $w_j = entmax_{\nu} \left(\frac{\theta_F}{T}\right)_j$.

the weights $\{w_1, ..., w_p\}$ to be trained as 0 or 1, but these weights may not all be 0 and 1. In other words, in NODE-GA¹M, NODTs may estimate the higher-order components rather than main effects. Therefore, it is difficult to consider NODE-GAM as the functional ANOVA model which decomposes a high-dimensional function into the sum of low-dimensional functions. Furthermore, in NODE-GAM, we can not obtain the estimated component function.

However, the feature function F^c of ANOVA-NODT for component S at depth c uses only the input features corresponding to S. For c = 1, ..., |S|, we use feature function defined as

$$F^c(\mathbf{x}) = (\mathbf{x}_S)_c$$

2371 2372 2373

2370

2374 ANOVA-NODE estimates component f_S by an ensemble of ANOVA-NODTs corresponding to S. 2375 Therefore, ANOVA-NODE is a functional ANOVA model. Therefore, although both NODE-GAM and ANOVA-NODE utilize NODT, they are fundamentally different models.







Figure 30: Plots of the functional relations of the main effects on 5 randomly sampled training data from WINE datasets.

2430 M ON THE POST-PROCESSING FOR THE SUM-TO-ZERO CONDITION

2432 2433

2434

2435

We have seen that NA^2M and NB^2M are competitive in prediction performance even though they are poor in estimating the components. There is a way to transform any estimate of a component to one that satisfies the sum-to-zero condition (Lengerich et al., 2020).

We consider a estimated $GA^{d}M \hat{f}(\mathbf{x}) = \beta_0 + \sum_{k=1}^{d} \sum_{S_k \subseteq [p]} \hat{f}_{S_k}(\mathbf{x}_{S_k})$ where S_k is a component for $|S_k| = k$. We write f_S instead of $f_S(\mathbf{x}_S)$ for notational simplicity since \mathbf{x} is fixed. First of all, for component S_d , we can transform \hat{f}_{S_d} into

2440 2441

2444 2445

where \tilde{f}_{S_d} satisfies the sum-to-zero condition. Next, for k = 1, ..., d, we redefine $\hat{f}_{S_{d-k}}$ into

$$\hat{f}_{S_{d-k}} = \hat{f}_{S_{d-k}} - (-1)^{d-k} \int_{\mathcal{X}_{S_k}} \hat{f}_{S_d} d\Pi_{j \in S_k} \mu_j$$

 $\tilde{f}_{S_d} = \hat{f}_{S_d} + \sum_{k=1}^d \sum_{V \subseteq S_d, |V|=k} (-1)^{d-k} \int_{\mathcal{X}_V} \hat{f}_{S_d} d\Pi_{j \in V} \mu_j$

where $S_{d-k} = S_d \setminus S_k$. If this process is performed sequentially for all components in order, all \tilde{f}_{S_k} terms in $\hat{f}(\mathbf{x}) = \tilde{\beta}_0 + \sum_{k=1}^d \sum_{S_k \subseteq [p]} \tilde{f}_{S_k}(\mathbf{x}_{S_k})$ satisfy the sum-to-zero condition.

Let us consider performing post-processing for GA^dM on a given dataset. The computational order for post-processing of a single point **x** is $\mathcal{O}(dn^{d-1})$. Therefore, if post-processing is carried out for all data points, the computational order becomes $\mathcal{O}(dn^d)$. In other words, not only global interpretation (e.g., *l*1 norm, functional relation plots, etc.) but also local interpretation is practically infeasible. Furthermore, performing post-processing requires storing a dataset, which causes memory efficiency issues.

Additionally, in GA^dM, when post-processing is performed, and then calculating the ANOVA-SHAP $\phi_j(\mathbf{x})$ for a given point x requires a computational order of $\mathcal{O}(p^{d-1}dn^d)$, which is even more demanding.

Table 27 compares the stability scores of the main effects of 'Latitude' and 'Longitude' for ANOVA-N²ODE, NA²M and NB²M, and Figure 31 draws the 5 functional relations of the main effects of 'Latitude' and 'Longitude' estimated by ANOVA-N²ODE, NA²M and NB²M on 5 randomly sampled training data. It is observed that NA²M and NB²M are still unstable even after the post-processing, which indicates that instability in NAM and NBM is not only from unidentifiability but also instability of DNN.

The situation becomes different when we apply the post-processing to GAM-NODE. Table 28 presents the stability scores of ANOVA-N²ODE and post-processed GAM-N²ODE normalized by the stability score of ANOVA-N²ODE, and Figure 32 presents the plots of the functional relations of the main effects estimated by GAM-N²ODE on 5 randomly sampled training data. Interestingly, unlike NAM and NBM, it is observed that GAM-N²ODE becomes more stable after post-processing.

The post-processing would not be a practically usable method since computation cost is too large for calculating the integration. The order of computation for the post-processing GAM-N^dODE is $O(n^d)$, where n is the number of data points to which the post-processing is applied.

2479

2480 Table 27: Stability scores for 'Latitude' and 'Longitude'. 2481 ANOVA-N²ODE NA²M NB²M Model 2482 0.104 Latitude 0.006 0.067 2483 0.015 Longitude 0.094 0.103



Figure 32: Plots of the functional relations of the main effects in post-processed GAM-N²ODE on CAL-HOUSING dataset.

2538 N DISCUSSION OF RELATED LITERATURE.2539

In this section, we describes the comparison between ANOVA-NODE and additive higher-order
 factorization machines (AHOFMs, Rügamer (2024)).

Higher-Order Factorization Machines (HOFMs). Higher-Order Factorization Machines (HOFMs, Blondel et al. (2016)) are a Factorization Machine(FM) which considers up to higher-order interaction. Blondel et al. (2016) proposed an algorithm that efficiently learns HOFMs using the properties of the ANOVA kernel, even as the input dimension p increases.

Additive higher-order factorization machines (AHOFMs). AHOFMs (Rügamer, 2024)) approximate each component in the functional ANOVA model using a spline basis representation. That is, for j = 1, ..., p, AHOFMs estimate f_j by

2551 2552

2553

2547

 $f_j(x_j) = \sum_{m=1}^{M_j} B_{m,j}(x_j) \beta_{m,j}$

where $\{B_{1,j}, ..., B_{M_j,j}\}$ are the basis functions and $\{b_{1,j}, ..., b_{M_j,j}\}$ are the coefficients. Furthmore, for high-order interaction *S*, AHOFM estimates f_S via extended spline basis representation which used the tensor product spline. Finally, Rügamer (2024) proposed an algorithm for efficiently learning AHOFMs, which applies the method used in HOFMs to AHOFMs, taking higher-order interactions into account.

2560
 2561
 2561
 2562
 2562
 2564
 2564
 2565
 2565
 2565
 2566
 2566
 2567
 2568
 2568
 2569
 2569
 2569
 2569
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2561
 2562
 2562
 2562
 2562
 2562
 2563
 2564
 2564
 2564
 2564
 2564
 2565
 2565
 2564
 2564
 2565
 2565
 2564
 2564
 2565
 2565
 2564
 2564
 2565
 2565
 2565
 2565
 2565
 2565
 2565
 2565
 2566
 2567
 2567
 2568
 2568
 2568
 2569
 2569
 2569
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 2560
 <li

If all higher-order interactions are considered, the method proposed by Rügamer (2024)) is more efficient in terms of scalability compared to ANOVA-NODE. However, the experiments in this paper show that considering interactions up to the second order is sufficient for real data. Moreover, when all higher-order interactions are considered in the functional ANOVA model, it is likely to be overfitted.

Since ANOVA-NODT is also a model represented by the spline basis representation of equation (2) in Rügamer (2024)(*entmax* can be viewed as basis function) factorization method can be used to improve the scalability of ANOVA-NODE. Applying the method used by AHOFMs to ANOVA-NODE seems like an interesting research topic, but in this paper, we proposed NBM-NODE in Appendix G.2, which enhances the scalability of ANOVA-NODE by utilizing the method by which Radenovic et al. (2022) improved the scalability of NAM (Agarwal et al. (2021)).

There are two main differences between ANOVA-NODE and AHOFMs: one is whether the basis function is learned, and the other is whether the sum-to-zero condition is satisfied. To be specific, AHOFMs do not learn the basis function, but ANOVA-NODE can be seen as learning the basis function by learning the parameters γ and b in the *entmax* function. Also, since AHOFMs do not satisfy the sum-to-zero condition, the estimated components are not identifiable. However, unlike AHOFMs, the estimated components by ANOVa-NODE are identifiable. That is, ANOVA-NODE provides more reliable interpretations than AHOFMs.

- 2581 2582
- 2583
- 2584
- 2585
- 2586
- 2587
- 2588
- 2589
- 2590