MIXNAM: ADVANCING NEURAL ADDITIVE MODELS WITH MIXTURE OF EXPERTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Additive models, such as Neural Additive Models (NAMs), are recognized for their transparency, providing clear insights into the impact of individual features on outcomes. However, they traditionally rely on point estimations and are constrained by their additive nature, limiting their ability to capture the complexity and variability inherent in real-world data. This variability often presents as different influences from the same feature value in various samples, adding complexity to prediction models. To address these limitations, we introduce MixNAM, an innovative framework that enriches NAMs by integrating a mixture of experts, where each expert encodes a different aspect of this variability in predictions from each feature. This integration allows MixNAM to capture the variability in feature contributions through comprehensive distribution estimations and to include feature interactions during expert routing, thus significantly boosting performance. Our empirical evaluation demonstrates that MixNAM surpasses traditional additive models in performance and is comparable to complex black-box approaches. Additionally, it improves the depth and comprehensiveness of feature attribution, setting a new benchmark for balancing interpretability with performance in machine learning. Moreover, the flexibility in MixNAM configuration facilitates the navigation of its trade-offs between accuracy and interpretability, enhancing adaptability to various data scenarios.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

032 Deep neural networks (DNNs) have proven exceptional at modeling complex data relationships, 033 achieving remarkable results across various domains, such as computer vision and natural language 034 processing (Young et al., 2018; Hassaballah & Awad, 2020). Despite these successes, the opaque 035 nature of DNN outputs often makes them non-interpretable, which hinders their usage in high-stakes areas like healthcare and finance where the basis of decisions is crucial (Ghnemat et al., 2023; Liao 037 & Varshney, 2021). This critical limitation has spurred research into interpretable machine-learning 038 models (Alvarez Melis & Jaakkola, 2018; Ghorbani et al., 2019; Wu et al., 2020; Yeh et al., 2020). These models aim to make the reasoning behind decisions more transparent, facilitating trust and wide adoption in sensitive applications. 040

Among the interpretable approaches developed, additive models stand out for their transparency, enabling direct visualization of how individual features impact predictions (Hastie, 2017). Additive models independently encode each feature through a non-linear transformation and subsequently combine them linearly for final predictions, simplifying the understanding of the influence from each feature. Based on the framework of additive models, Neural Additive Models (NAMs; Agarwal et al., 2021) attempt to model the complex patterns of each feature with DNNs, offering a flexible and differentiable representation of features while maintaining the additive structure.

While the additive nature of feature information in NAMs allows for clear explanations of individ ual features, this structure inherently limits their ability to learn intricate mappings between input
 features and target outputs. Thus, the performance of these models often lags behind that of their
 complex black-box counterparts, which can learn interactions among features. Furthermore, addi tive models traditionally produce a single outcome for each feature value, an approach that fails to
 account for the variability in outcomes that occur in complex real-world scenarios. The variability
 model prediction refers to the extent to which the predicted outputs of a model fluctuate or differ

when the same input feature value is presented across different samples. For example, the same health metric might have different implications depending on an individual's age group or other contextual factors. By relying on point-to-point estimations, these models do not capture the actual distribution of possible outcomes associated with a given feature value. Recent attempts to enhance NAMs with uncertainty estimation seek to address this issue (Bouchiat et al., 2023; Thielmann et al., 2024). However, these approaches are still limited by their inevitable assumptions of prior distributions and the additive nature of the entire models, which restrict their ability to accurately reflect complex underlying distributions and interactions.

062 To address the limitations of traditional additive models while retaining their interpretative bene-063 fits, we introduce MixNAM, a novel and general framework that extends beyond the constraints of 064 NAMs. MixNAM integrates multiple expert feature encoders operating in parallel, each of which encodes a different aspect of the variability in model predictions. Moreover, a dynamic routing 065 mechanism is introduced to assess and combine the relevance of different experts for the final pre-066 diction. Our experiments demonstrate that MixNAM significantly outperforms all existing additive 067 models, matching the performance of complex black-box methods. Different from traditional addi-068 tive models that offer point estimations, MixNAM provides a more nuanced understanding of feature 069 impacts by visualizing the range and distribution of possible outcomes for each feature value. The results further reveal that MixNAM can be seen as a generalization of NAM, with the flexibility to 071 adjust the balance between detailed distribution estimations and straightforward point estimations, catering to specific needs for transparency and performance. The contributions of this work can be 073 summarized as follows: 074

- We propose MixNAM, a general framework built upon NAMs, by incorporating a mixture of experts with dynamic routing algorithms.
- MixNAM demonstrates superior performance compared to traditional additive models and is comparable to complex black-box models.
- MixNAM preserves the capability of explaining feature contributions through visualizations, enhancing the current point estimations with nuanced distribution estimations.
- MixNAM enables a seamless transition from detailed distribution estimation to precise point estimation, effectively balancing model accuracy and interpretability.

2 BACKGROUND

2.1 GENERALIZED ADDITIVE MODELS

Additive models, or Generalized Additive Models (GAMs; Hastie, 2017), represent an evolving area in the field of explainable artificial intelligence, blending interpretability with predictive accuracy. Given a sample with n features x_1, \dots, x_n , the task of an additive model, in general, is to learn corresponding feature encoders f_1, \dots, f_n which map the input values from feature domains to the prediction domain. The final predicted value of the target output y is computed as

$$\hat{y} = w_0 + \sum_{i=1}^n f_i(x_i), \tag{1}$$

where w_0 is the learnable bias of the output distribution. GAMs highlight the influence of individual features additively, facilitating easy comprehension and interpretation. Advances such as Explainable Boosting Machine (EBM; Lou et al., 2012) and Neural Oblivious Decision Ensembles for GAM (NODE-GAM; Chang et al., 2022) have built upon this foundation by integrating boosting techniques (Friedman, 2001) and differentiable decision trees (Popov et al., 2020), respectively, enhancing accuracy while preserving interpretability.

103 104

075

076 077

078

079

081

082

084

085

087

091

092

094 095

2.2 NEURAL ADDITIVE MODELS

Developed from traditional additive models, Neural Additive Models (NAMs; Agarwal et al., 2021)
 incorporate neural networks to encode each feature, marking a significant evolution in the modeling
 capability of GAMs. This neural integration allows NAMs to capture more complex feature relation ships without sacrificing the interpretative benefits. Subsequent developments, such as Neural Basis

108 Models (NBMs; Radenovic et al., 2022) and Gaussian Process Neural Additive Models (GP-NAMs; 109 Zhang et al., 2024) have further refined this approach. They focus on efficiency and scalability, re-110 ducing the complexity and enhancing the ability of models to manage larger datasets with complex 111 underlying structures. Recent research has also explored extending NAMs to include uncertainty 112 estimation, providing a more comprehensive explanation of feature influence by considering distribution parameters and Bayesian inference methods (Bouchiat et al., 2023; Thielmann et al., 2024). 113 However, these models generally rely on a simplistic assumption about output distributions, which 114 may limit their effectiveness in real-world scenarios where data interactions are complex. 115

116

117 2.3 MIXTURE OF EXPERTS

118 The Mixture of Experts (MoE) technique has gained substantial attention from the research com-119 munity recently for its ability to enhance both the scalability and expertise of neural networks. It is 120 designed to learn multiple specialized sub-models, known as "experts", on different subtasks. A gat-121 ing network is trained to dynamically assign a new input to the most relevant experts (Jacobs et al., 122 1991). One key advantage of the MoE model is its efficiency and adaptability in handling large-scale 123 problems. Shazeer et al. (2017) explored scaling MoE for applications in massive neural network 124 architectures by incorporating a sparse gating network to select only a small subset of experts for 125 each input. They demonstrate that integrating MoE layers can significantly boost model capacity with a minor loss in computational efficiency. Recent work also explores the adaptation of MoE in 126 the Transformer architecture (Vaswani et al., 2017), encouraging their use in cutting-edge research 127 such as large language models (Artetxe et al., 2022; Jiang et al., 2024; Lepikhin et al., 2021) and 128 vision language models (Lin et al., 2024; Shen et al., 2023). 129

- 130
- 131 132

152 153

154

3 MIXNAM: MIXTURE OF NEURAL ADDITIVE MODELS

Our Mixture of Neural Additive Models (MixNAM) framework extends traditional Neural Additive Models (NAMs) by incorporating multiple expert predictors for each feature to better capture the variability and complexity of data relationships. This approach combines feature-specific expert encoders with a dynamic expert routing mechanism to determine the most relevant experts for each input scenario, as depicted in Figure 1.





3.1 FEATURE ENCODING WITH EXPERTS

In traditional additive models, a single predictor associated with each feature often limits the expressiveness and adaptability of the model in complex scenarios. MixNAM addresses this by embedding a mixture of C expert predictors for each feature x_i ($i \in \{1, \dots, n\}$), significantly enhancing the model's capacity to capture diverse non-linear relationships within the data. Specifically, each feature x_i is encoded by an expressive feature encoder f_i , which transforms the input into a latent vector. We then enhance the standard encoding process by introducing C expert predictors E_{i1}, \dots, E_{iC} , where each predicts different potential outcomes from the encoded information of x_i :

$$o_{ik} = E_{ik}(f_i(x_i)), \quad \forall i \in \{1, \cdots, n\}, k \in \{1, \cdots, C\}.$$
 (2)

This architecture allows each feature to express a spectrum of influences, capturing subtle variations in the data that might be overlooked by more monolithic approaches. In our experiments, each f_i is implemented as a multi-layer neural network, and E_{ik} is implemented as a linear layer.

166 3.2 DYNAMIC EXPERT ROUTING

171 172

173

177

190

191 192

201

210 211 212

Expert routing is pivotal in MixNAM, where the relevance of each expert is dynamically assessed through a routing mechanism. MixNAM uses a router to compute relevance scores for experts based on the overall context of input and aggregates expert predictions for final outputs:

$$o_i = \sum_{k=1}^{C} r_{ik} o_{ik}$$
 w.r.t. $r_{ik} \ge 0$ and $\sum_{k=1}^{C} r_{ik} = 1$, $\forall i \in \{1, \cdots, n\}$, (3)

where r_{ik} is the estimated relevance for the k-th expert in feature x_i , and o_i is the summarized outcome given by x_i . The final prediction \hat{y} is generated by adding outputs o_1, o_2, \dots, o_n from each feature with an additional bias term ω_0 similar to Formula 1.

178 Score Estimation. Motivated by the fact that one spe-179 cific feature input may correspond to different target values depending on the overall information of the instance, we propose a routing mechanism that takes all features as 181 input and estimates the relevance of each expert for pre-182 diction. Suppose the encoded information $f_i(x_i)$ from f_i 183 is represented by a d-dimensional vector. For the C experts of the feature x_i , the routing system in MixNAM 185 learns a $d \times C$ scoring matrix \mathcal{A}_{ij} to assess how infor-186 mation encoded from x_i affects the expert selection of x_i . 187 The summarized scores for expert relevance estimation in 188 feature x_i can be produced by aggregating the information 189 from all features:

$$\varphi_j = \mu_j + \sum_{i=1}^n \mathcal{A}_{ij}^\top f_i(x_i) \tag{4}$$



Figure 2: An illustration of our routing mechanism. This example contains two features with three experts each.

(6)

193 where μ_j is a learnable *C*-dimensional bias term. The 194 scores in φ_j are then normalized into a valid relevance dis-195 tribution over experts using an activation strategy.

As depicted in Figure 2, all scoring matrices can be organized as a large block matrix, where the block at the *i*-th row and the *j*-th column corresponds to A_{ij} . Leveraging parallel computing resources, such as GPUs, this organization facilitates accelerated score estimation by processing multiple features parallelly. We also explore a specialized version of MixNAM in Appendix G, where the scoring matrices compose a block diagonal matrix, making it a strictly additive model.

Activation Strategy. To effectively transform the raw scores into a valid distribution of expert relevance, an activation strategy is needed to project $\varphi_i \in \mathbb{R}^C$ to the simplex $\mathcal{D} = \{v \in \mathbb{R}^C_+ | v^\top 1 = 1\}$. This ensures that the relevance scores fulfill the constraints required by the model, where each score vector must sum to one and contain only non-negative values.

We utilize a sparse activation strategy adapted from previous MoE research (Jiang et al., 2024) to focus on the most significant experts without overwhelming the model with redundant information. Given scoring vector φ_i , our model will first create a masking vector $M_i \in \mathbb{R}^C$ to zero out less significant scores, whose k-th entry is defined as

$$M_i[k] = \begin{cases} 0, & \text{if } \varphi_i[k] \text{ is among the largest } K \text{ values in } \varphi_i, \\ -\infty, & \text{otherwise,} \end{cases}$$
(5)

where K is the number of experts activated. The relevance scores will then be computed by

214
215
$$r_{ik} = \exp(\varphi_i[k] + M_i[k]) / \sum_{l=1}^{C} \exp(\varphi_i[l] + M_i[l]).$$

Such an activation strategy results in a continuous subset of values in \mathcal{D} , allowing for a nuanced learning of expert relevance. While the router is flexible for learning a continuous range of outputs, our analysis in Appendix F proves that it is not a universal function approximator but is inherently a generalized additive model with normalization. In Appendix H, we also discuss an alternative activation strategy that brings a finite set of points in \mathcal{D} as a discrete range for relevance activation. It shows distinct model behaviors that the MixNAM framework achieves with different configurations.

222 223

224

236

237 238

239

240241242243244

245

246

247 248

249

3.3 Optimization and Training

The training of MixNAM involves optimizing standard loss functions adapted to downstream tasks, 225 such as cross-entropy or mean squared error. Following previous studies (Agarwal et al., 2021; 226 Radenovic et al., 2022), we employ the dropout and L2 regularization on model parameters to avoid 227 overfitting and promote robust generalization. An L2 penalty of feature outcomes is also applied 228 to suppress redundancy. Additionally, we introduce an expert dropout to prevent the model from 229 overly depending on specific routing paths. An expert variation penalty is also introduced, designed 230 to minimize variance in predictions across experts for the same input feature, encouraging the model 231 to learn meaningful and interpretable expert functions. This penalty is the key to balance accuracy 232 and interpretability in MixNAM, which is discussed in detail in Section 4.3.

As introduced in Section 3.2, the final prediction of an input sample by MixNAM is obtained by linearly aggregating the predicted outcomes from each feature

$$\hat{y} = \omega_0 + \sum_{i=1}^n o_i = \omega_0 + \sum_{i=1}^n \sum_{k=1}^C r_{ik} o_{ik}.$$
(7)

Suppose we have N training samples x^1, \dots, x^N with corresponding task labels y^1, \dots, y^N . The objective during training is

minimize
$$\frac{1}{N} \sum_{t=1}^{N} \mathcal{L}(y^t, \hat{y}^t) + \frac{\gamma}{nN} \sum_{t=1}^{N} \sum_{i=1}^{n} (o_i^t)^2 + \frac{\lambda}{nNC} \sum_{t=1}^{N} \sum_{i=1}^{n} \sum_{k=1}^{C} \left(o_{ij} - \frac{o_{i1} + \dots + o_{iC}}{C} \right)^2,$$
(8)

where \hat{y}^t denotes the final model prediction for x^t ($t \in \{1, \dots, N\}$). γ is the weight for the output penalty and λ is the weight for the expert variation penalty. The loss function \mathcal{L} can be the mean square error loss or the cross-entropy loss, depending on the task that the model is dealing with.

4 EXPERIMENTS

250 The experiments in this section are designed to comprehensively evaluate MixNAM across dimen-251 sions of accuracy and interpretability, core attributes that define its utility in practical applications. 252 We begin by assessing the accuracy of MixNAM through performance comparisons with established 253 baselines ranging from traditional additive models to complex black-box models. We then delve into 254 the interpretability aspect, utilizing visualizations to demonstrate how MixNAM elucidates the im-255 pact of individual features on predictions. Moreover, experiments are conducted to showcase the 256 unique ability of MixNAM to balance accuracy with interpretability, illustrating how this can be ad-257 justed to meet specific application demands. In addition to the experiments on real-world datasets, we test our MixNAM with simulated data to see if the MoE design helps it capture multimodal data. 258 Additional results in Appendix I reveal how the performance of MixNAM scales with the number 259 of experts in its configuration. 260

261 262

263

264

4.1 EVALUATION OF MIXNAM ACCURACY

4.1.1 DATASETS

To evaluate the accuracy of MixNAM, we select six widely recognized datasets that span a mix of
regression and classification tasks: Housing (Pace & Barry, 1997), MIMIC-II (Saeed et al., 2011),
MIMIC-III (Johnson et al., 2016), Income (Blake, 1998), Credit, and Year. These datasets are chosen
due to their varying complexities, including differences in the number of instances, features, and the
presence of categorical variables, ensuring a comprehensive assessment of the model performance
across diverse data scenarios. Detailed information about each dataset can be found in Appendix A.

270 4.1.2 BASELINES

272 To evaluate the accuracy of MixNAM, we benchmark it against a range of models selected for their relevance to the goals of MixNAM. This includes linear and spline models, which provide basic fea-273 ture explanation capabilities similar to additive models, Neural Additive Models (NAMs; Agarwal 274 et al., 2021) and Neural Basis Models (NBMs; Radenovic et al., 2022), which use neural networks 275 to enhance feature encoding, Explainable Boosting Machine (EBM; Lou et al., 2012; Nori et al., 276 2019) and NODE-GAM (Chang et al., 2022), which integrate advanced techniques like gradient 277 boosting (Friedman, 2001) and differentiable decision trees (Popov et al., 2020). Additionally, ex-278 tended additive models such as EB²M, NA²M, NB²M, and NODE-GA²M, which encode pairwise 279 feature interactions, are assessed. While these models facilitate interaction-based explanations, they 280 lack the capacity to succinctly explain the impact of individual features, presenting challenges in 281 attribute summarization due to potential dimensionality issues. Traditional black-box models like 282 MLP, NODE (Popov et al., 2020), and XGBoost (Chen & Guestrin, 2016) are also evaluated for 283 their ability to flexibly encode feature interactions, though at the cost of interpretability. The imple-284 mentation details of our experiments are presented in Appendix B and C.

4.1.3 EVALUATION RESULTS

Table 1 presents the comparison of MixNAM against all baseline models. Evaluation metrics are provided below each dataset name, with arrows indicating the desired direction for scores. For each dataset, the top four scores are underlined to highlight a broad range of competitive models. In addition to the evaluation scores, a column labeled "FA" (Feature Attribution) is added to the table, indicating whether the method can provide explanations focused on individual feature impacts.

293 294

295

296

285

Table 1: Performance comparison on benchmark datasets. "FA" (Feature Attribution) capabilities of different models are denoted by ✓ for presence and ✗ for absence. The top four scores for each dataset are underlined. MixNAM significantly outperforms traditional additive models with "FA".

		Housing	MIMIC-II	MIMIC-III	Income	Credit	Year
Model	FA	RMSE↓	AUC ↑	AUC ↑	AUC ↑	$\overline{\text{AUC}\uparrow}$	$\overline{\text{MSE}\downarrow}$
MLD	~	0.501	0.835	0.815	0.914	0.981	78.48
MLP	^	± 0.006	± 0.014	± 0.009	± 0.003	± 0.007	± 0.56
NODE	v	0.523	0.843	0.828	0.919	0.981	<u>76.21</u>
NODE	<u>^</u>	± 0.000	± 0.011	± 0.007	± 0.003	± 0.009	± 0.12
VGBoost	Y	<u>0.443</u>	0.844	0.819	<u>0.928</u>	0.978	<u>78.53</u>
AUDOOSI	<u>^</u>	± 0.000	± 0.012	± 0.004	± 0.003	± 0.009	± 0.09
ED^2M	~	0.492	0.848	0.821	0.928	0.982	83.16
ED IVI	<u>^</u>	± 0.000	± 0.012	± 0.004	± 0.003	± 0.006	± 0.01
$\mathbf{M} \mathbf{A}^{2} \mathbf{M}$	v	0.492	0.843	0.825	0.912	0.985	79.80
INA M	<u>^</u>	± 0.008	± 0.012	± 0.006	± 0.003	± 0.007	± 0.05
ND^2M	v	0.478	0.848	0.819	0.917	0.978	79.01
IND IVI	^	± 0.002	± 0.012	± 0.010	± 0.003	± 0.007	± 0.03
NODE $C \Lambda^2 M$	v	<u>0.476</u>	<u>0.846</u>	0.822	0.923	<u>0.986</u>	79.57
NODE-GA M	^	± 0.007	± 0.011	± 0.007	± 0.003	± 0.010	± 0.12
Lincor		0.735	0.796	0.772	0.900	0.976	88.89
Lineal		± 0.000	± 0.012	± 0.009	± 0.002	± 0.012	± 0.40
Spline		0.568	0.825	0.812	0.918	0.982	85.96
Spine		± 0.000	± 0.011	± 0.004	± 0.003	± 0.011	± 0.07
FRM		0.559	0.835	0.809	0.927	0.974	85.81
EDIVI		± 0.000	± 0.011	± 0.004	± 0.003	± 0.009	± 0.11
NAM		0.572	0.834	0.813	0.910	0.977	85.25
INAIVI		± 0.005	± 0.013	± 0.003	± 0.003	± 0.015	± 0.01
NBM		0.564	0.833	0.806	0.918	0.981	85.10
INDIVI		± 0.001	± 0.013	± 0.003	± 0.003	± 0.007	± 0.01
NODE GAM		0.558	0.832	0.814	<u>0.927</u>	0.981	85.09
NODE-GAM		± 0.003	± 0.011	± 0.005	± 0.003	± 0.011	± 0.01
MiwNAM		<u>0.451</u>	<u>0.847</u>	<u>0.825</u>	<u>0.927</u>	<u>0.982</u>	<u>78.66</u>
IVIIALV/AIVI		± 0.002	± 0.014	± 0.006	± 0.003	± 0.009	± 0.21

From the table, it can be observed that traditional additive models, while capable of providing feature attributions, do not generally perform as well as models that can encode feature interactions, such as NA²Ms and various black-box models. However, MixNAM distinguishes itself by offering competitive performance, achieving comparable results to models encoding feature interactions, while also providing the valuable capability of feature-level explanation as indicated by the "FA" column. This notable combination of high performance and strong interpretability in MixNAM underscores its effectiveness in complex real-world applications.

4.2 INTERPRETABILITY OF MIXNAM PREDICTION

Beyond its superior performance relative to additive models, MixNAM can explain how individual features influence the final prediction as well. MixNAM utilizes a dynamic routing mechanism where the relevance of different experts for a feature is confined within a simplex. This constraint allows MixNAM to determine precise upper and lower prediction bounds for any given feature value x_i with

$$u_{i} = \max_{k \in \{1, \cdots, C\}} E_{ik}(f_{i}(x_{i})) \quad \text{and} \quad l_{i} = \min_{k \in \{1, \cdots, C\}} E_{ik}(f_{i}(x_{i})).$$
(9)

The bounds represent the maximum and minimum potential outputs for a feature, which ensure that the outputs remain within a definite and interpretable range. Additionally, by iterating over all instances, we can plot the actual feature influences after expert routing, showing how predictions distribute within the bounds.

344 Figure 3 illustrates the interpretative capability of MixNAM, using the "Longitude" feature from the 345 Housing dataset as an example. For comparison, we include visual explanations by EBM and NAM on this feature. We also include the "Longitude-Latitude" interaction plot from NA²M to validate 346 whether the distribution of predictions from MixNAM aligns with models that encode complex 347 interactions. In the figure of MixNAM, the upper bound is highlighted in dark blue, and the lower 348 bound is in red. The actual predicted values are plotted as blue dots, showing the entire distribution 349 of scores in the dataset. The y-axis in the figures of EBM, NAM, and MixNAM shows the mean-350 centered contribution to the model prediction given by different feature values. We also plot the 351 color bars in the background to reflect the normalized data density following prior studies (Agarwal 352 et al., 2021; Radenovic et al., 2022). In the NA²M figure, the contribution given by each pair of 353 feature values is reflected by the color of the corresponding dot. 354



359 360 361 362 363 364

365 366

355

356

357

358

331 332

333

339

Figure 3: Visual comparison of the "Longitude" feature impact on house prices across models.

The visualizations in Figure 3 reveal that EBM, NAM, and MixNAM all capture similar geographic 367 trends affecting house prices, with sharp increases observed around the key longitude markers at 368 -122.5 (122.5°W, San Francisco) and -118.5 (118.5°W, Los Angeles). Compared with EBM and 369 NAM, MixNAM provides a more comprehensive insight into the output distribution by leveraging 370 dynamic expert routing that incorporates all features. Specifically, while high prediction scores in 371 regions like San Francisco and Los Angeles align with other models, MixNAM captures a broad 372 range of possible outcomes between -120 and -119, with most data points clustering around the 373 lower bound. This detailed variance is validated by the "Longitude-Latitude" interaction analysis 374 from NA^2M , where coastal regions are shown to positively impact house prices (blue dots) while 375 other areas in the same longitude range generally lower property values (red dots). In contrast to NA²M, which requires the selection of specific interactions to convey such insights, MixNAM seam-376 lessly integrates and displays the complex underlying data distribution in a single plot, enhancing 377 interpretative clarity and effectiveness.

386

387

391

392

417 418

419

421

426 427

430 431

378 4.3BALANCE BETWEEN ACCURACY AND INTERPRETABILITY 379

380 In addition to demonstrating the accuracy and interpretability of MixNAM, we further explore how these two characteristics can be adjusted to reach a balance. As mentioned in Section 3.3, an expert 381 variation penalty is implemented during training to ensure that MixNAM develops meaningful ex-382 pert functions. Our qualitative and quantitative analyses demonstrate that, by adjusting the penalty 383 weight λ , MixNAM presents a general framework of "constrained interpretable models", which 384 provides an opportunity to control the trade-offs between model accuracy and its interpretability. 385

4.3.1 QUALITATIVE ANALYSIS OF SHAPE PLOTS BY MIXNAM

388 While the default value of the variation penalty weight λ is 0.1, we rerun our experiments on the 389 Housing dataset with diverse weight values from 0 to 100, examining how the model performance 390 and explanation change with different λ . Table 2 illustrates the impact of varying expert variation penalties on MixNAM. It reveals that MixNAM remains stable performance when $\lambda < 0.1$. However, there is a noticeable decline in model performance as λ increases beyond the threshold.



Consistent with the stable performance for $\lambda \leq 0.1$, it is shown by the shape plots in Table 2 407 that the distribution of the predicted values remains similar across such settings. However, as λ 408 decreases further, the gap between the upper and lower bounds widens, resulting in looser bounds 409 that become less informative and fail to accurately reflect the true limits of the data distribution. 410 Conversely, larger λ values tighten these bounds, better reflecting the distribution that the model 411 captures. Nevertheless, excessively high λ values may also constrain the ability of MixNAM to 412 learn complex underlying data distribution, shifting the focus from estimating distributions to point 413 estimation. The learned shape plot becomes almost identical to the ones generated by traditional 414 additive models when the variation penalty weight is set too large (e.g., $\lambda = 100$). The varying shape 415 plots qualitatively show how the variation penalty controls the trade-off between model accuracy and 416 interpretability, providing different visualization results tailored for various interpretation needs.

4.3.2 QUANTITATIVE ANALYSIS OF MODEL ADDITIVITY AND BOUND TIGHTNESS

In addition to the qualitative analysis of shape plots, we design two metrics, model additivity and 420 bound tightness, to quantitatively explain how the trade-off can be controlled by the variation penalty in MixNAM. To show how much additivity is preserved in MixNAM and how close it is to a strict 422 additive model, we look into the gradient of the model prediction with respect to each feature and 423 decompose its variation to examine the driven factors in the variation. The total variance of the 424 derivative of \hat{y} with respect to feature x_i can be decomposed as 425

$$\operatorname{Var}\left(\frac{\partial \hat{y}}{\partial x_i}\right) = \mathbb{E}\left(\operatorname{Var}\left(\frac{\partial \hat{y}}{\partial x_i}|x_i\right)\right) + \operatorname{Var}\left(\mathbb{E}\left(\frac{\partial \hat{y}}{\partial x_i}|x_i\right)\right).$$
(10)

428 The additivity of a model on a feature is measured by quantifying how much of its derivative is 429 determined by the feature itself, and the additivity of the entire model can be computed as:

$$\operatorname{Additivity} = \frac{1}{n} \sum_{i=1}^{n} \frac{\operatorname{Var}(\mathbb{E}(\frac{\partial \hat{y}}{\partial x_i} | x_i))}{\operatorname{Var}(\frac{\partial \hat{y}}{\partial x_i})}, \tag{11}$$

432 where x_1, \dots, x_n are the *n* features considered. The value of an additivity score will range from 0 433 to 1, where 1 means strictly additive and 0 means non-additive. In addition to the model additivity, 434 we also examined the tightness of estimated bounds given by MixNAM in the shape plots, which is 435 an important factor when interpreting the results. The bound tightness is measured by

$$\text{Tightness} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E} \Big[\frac{\max(o_i | x_i) - \min(o_i | x_i)}{\operatorname{upper}(o_i | x_i) - \operatorname{lower}(o_i | x_i)} \Big],$$
(12)

where the o_i is the predicted output given by the feature encoders of x_i as defined in Formula 3. 440 This metric measures how tightly the bounds fit actual model outputs. 441

Table 3 presents the model additivity and bound tightness of MixNAM with different λ values. As the additivity metric is generalizable to all differentiable models, we include the measured additivity of MLP, NA^2M , and NAM for a comprehensive comparison. The results show a clear trade-off between model additivity and performance, as models with high additivity values tend to perform worse on the given task. By mitigating the additivity restrictions, MixNAM has significantly better performance with reduced additivity, which, nevertheless, is still much higher than models without feature attributions such as MLP and NA²M. Moreover, the results quantitatively show that λ can adjust the trade-off between additivity and accuracy, with model additivity close to 1 when λ is large.

Table 3: Quantitative analysis of model additivity and bound tightness on the Housing dataset.

	MLP	NA ² M	$\begin{array}{c} \text{MixNAM} \\ (\lambda = 0) \end{array}$	$\begin{array}{c} \text{MixNAM} \\ (\lambda = 0.1) \end{array}$	$\begin{array}{c} \text{MixNAM} \\ (\lambda = 1) \end{array}$	$\begin{array}{c} \text{MixNAM} \\ (\lambda = 10) \end{array}$	$\begin{array}{l} \text{MixNAM} \\ (\lambda = 100) \end{array}$	NAM
Additivity	0.085	0.170	0.366	0.381	0.416	0.746	0.998	1.000
Additivity	± 0.011	± 0.020	± 0.003	± 0.013	± 0.003	± 0.003	± 0.000	± 0.000
Tightness			0.860	0.953	0.988	0.999	1.000	
rightness	-	-	± 0.000	± 0.011	± 0.000	± 0.000	± 0.000	_
DMCE (1)	0.501	0.492	0.451	0.451	0.458	0.515	0.582	0.572
KNISE (\downarrow)	± 0.006	± 0.008	± 0.003	± 0.003	± 0.003	± 0.008	± 0.008	± 0.005

Similarly, the table shows that the tightness score in MixNAM grows with the increase of the variation penalty, which already achieves 0.953 when $\lambda = 0.1$. However, a lower tightness score does not necessarily correspond to better performance, as the bounds may be too loose to reflect the actual distribution of model outputs, affecting the interpretability of shape plots given by the estimated bounds. With the bound tightness metric, we can conveniently evaluate the interpretation quality of MixNAM given by its upper and lower bounds without plotting the actual shape functions.

SIMULATION STUDY ON UNIMODAL AND MULTIMODAL DATA 4.4

469 The Mixture of Experts (MoEs) is known to be effective in capturing multi-distribution or mul-470 timodal data, which can perform better than unimodal models in certain complex data scenarios 471 (Shazeer et al., 2017). To examine whether the MoE design in MixNAM helps it learn complex 472 multimodal data, we perform simulation studies on MixNAM and NAM using synthetic data with 473 known ground truth to test if MixNAM can capture the underlying patterns.

474 Suppose we have two random variables $x_1 \sim U(0,1)$ and $x_2 \sim Bernoulli(0.5) \times 2 - 1$. We 475 first simulate 10k samples following a simple unimodal data distribution, where the target y value is 476 generated as 477

$$y = \sin(4\pi x_1) + x_2 + \varepsilon$$
 with $\varepsilon \sim \mathcal{N}(0, 0.1^2)$ (13)

478 For the multimodal data, we generate 10k samples with two modes and define y as 479

$$y = \begin{cases} \sin(4\pi x_1) + x_2 + \varepsilon & \text{if} \quad (x_1, x_2) \in \{(x_1, x_2) | x_2 = 1\}, \\ -\sin(4\pi x_1) + x_2 + \varepsilon & \text{if} \quad (x_1, x_2) \in \{(x_1, x_2) | x_2 = -1\}. \end{cases}$$
(14)

481 482

480

442

443

444

445

446

447

448

449 450 451

462

463

464

465

466 467

468

Figure 4 shows the shape functions learned from the unimodal and multimodal data by NAM and 483 MixNAM, respectively. For the unimodal data where the data-generating process is additive, both 484 NAM and MixNAM effectively capture the relations between the features and the target output. 485 However, on the multimodal data where the relation between x_1 and y exhibits two distinct modes,

486 NAM fails to capture the complex data patterns and instead learns only a weak influence of the 487 feature on the output. Instead, our MixNAM effectively captures the intricate relation between 488 x_1 and y, predicting a periodic pattern in its upper and lower bounds, which aligns with the true 489 underlying data-generating process.



Figure 4: Shape plots of MixNAM and NAM on the simulated data.

In addition to the shape plots with estimated bounds given by MixNAM, we also visualize the 509 individual experts learned by MixNAM in Figure 5 to see if the expert functions capture the designed modes in our simulated data. For both data distributions, we train a MixNAM model with two experts both of which can be activated. The figure shows that both learned experts on the unimodal 512 data fit the actual data distribution and overlap with each other. It demonstrates that MixNAM 513 can behave as a strictly additive model if the data-generating process is truly additive. On the 514 multimodal data, the two experts in MixNAM capture similar periodic patterns but with opposite 515 values, aligning with the two modes in the true data distribution. Our simulation study demonstrates 516 that by incorporating MoE design into additive models, MixNAM effectively captures the complex patterns in multi-distribution and multimodal data, which the original NAM can hardly handle.



Figure 5: Learned experts in MixNAM for the simulated data.

5 CONCLUSION

507 508

510

511

517

518 519

521

522

523

524

525

527

528 529

530 531

532 533

We introduce MixNAM, a general framework that advances additive models by incorporating a 534 mixture of experts. Our experiments show that MixNAM overcomes the inherent limitations of tra-535 ditional additive models, enhancing accuracy while preserving interpretability. MixNAM not only 536 provides detailed explanations of feature output distributions but also allows for a flexible balance 537 between accuracy and interpretability, adapting to various application needs. This framework ex-538 tends the utility of additive models, providing advanced performance and insightful interpretability, which are crucial for real-world machine-learning applications.

540 REFERENCES

547

562

- 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.
- David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining
 neural networks. *Advances in neural information processing systems*, 31, 2018.
- Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Mihaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, et al. Efficient large scale language modeling with mixtures of experts. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11699–11732, 2022.
- Catherine Blake. Uci repository of machine learning databases. http://www. ics. uci. edu/~
 mlearn/MLRepository. html, 1998.
- Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020.
- Kouroche Bouchiat, Alexander Immer, Hugo Yèche, Gunnar Rätsch, and Vincent Fortuin. Improv ing neural additive models with bayesian principles. 2023.
- Chun-Hao Chang, Sarah Tan, Ben Lengerich, Anna Goldenberg, and Rich Caruana. How interpretable and trustworthy are gams? In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 95–105, 2021.
- Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. NODE-GAM: neural generalized addi tive model for interpretable deep learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 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.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- James Enouen and Yan Liu. Sparse interaction additive networks via feature interaction detection
 and sparse selection. Advances in Neural Information Processing Systems, 35:13908–13920,
 2022.
- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pp. 1189–1232, 2001.
- Rawan Ghnemat, Sawsan Alodibat, and Qasem Abu Al-Haija. Explainable artificial intelligence (xai) for deep learning based medical imaging classification. *Journal of Imaging*, 9(9):177, 2023.
- Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based
 explanations. *Advances in neural information processing systems*, 32, 2019.
- Mahmoud Hassaballah and Ali Ismail Awad. Deep learning in computer vision: principles and applications. CRC Press, 2020.
- Trevor Hastie and Robert Tibshirani. Generalized additive models for medical research. *Statistical methods in medical research*, 4(3):187–196, 1995.
- Trevor J Hastie. Generalized additive models. In *Statistical models in S*, pp. 249–307. Routledge, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

- 594 Stefan Hegselmann, Thomas Volkert, Hendrik Ohlenburg, Antje Gottschalk, Martin Dugas, and 595 Christian Ertmer. An evaluation of the doctor-interpretability of generalized additive models with 596 interactions. In Machine Learning for Healthcare Conference, pp. 46–79. PMLR, 2020. 597 Farzali Izadi. Generalized additive models to capture the death rates in canada covid-19. In Mathe-598 matics of Public Health: Proceedings of the Seminar on the Mathematical Modelling of COVID-19, pp. 153-171. Springer, 2021. 600 601 Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of 602 local experts. Neural computation, 3(1):79-87, 1991. 603 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In 604 International Conference on Learning Representations, 2016. 605 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-607 ford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 608 Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024. 609 Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad 610 Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, 611 a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016. 612 613 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25, 2012. 614 615 Friedrich Leisch. Flexmix: A general framework for finite mixture models and latent glass regres-616 sion in r. 2004. 617 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, 618 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional 619 computation and automatic sharding. In 9th International Conference on Learning Representa-620 tions, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. 621 622 Q Vera Liao and Kush R Varshney. Human-centered explainable ai (xai): From algorithms to user 623 experiences. arXiv preprint arXiv:2110.10790, 2021. 624 Bin Lin, Zhenyu Tang, Yang Ye, Jiaxi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and 625 Li Yuan. Moe-llava: Mixture of experts for large vision-language models. arXiv preprint 626 arXiv:2401.15947, 2024. 627 628 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Confer-629 ence on Learning Representations, 2018. 630 Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligible models for classification and regression. 631 In Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and 632 data mining, pp. 150-158, 2012. 633 634 Harsha Nori, Samuel Jenkins, Paul Koch, and Rich Caruana. Interpretml: A unified framework for 635 machine learning interpretability. arXiv preprint arXiv:1909.09223, 2019. 636 R Kelley Pace and Ronald Barry. Sparse spatial autoregressions. Statistics & Probability Letters, 33 637 (3):291–297, 1997. 638 639 Eric J Pedersen, David L Miller, Gavin L Simpson, and Noam Ross. Hierarchical generalized additive models in ecology: an introduction with mgcv. PeerJ, 7:e6876, 2019. 640 641 Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for 642 deep learning on tabular data. In 8th International Conference on Learning Representations, 643 ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. 644 645 Filip Radenovic, Abhimanyu Dubey, and Dhruv Mahajan. Neural basis models for interpretability. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in 646
- In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in
 Neural Information Processing Systems, volume 35, pp. 8414–8426. Curran Associates, Inc., 2022.

648 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the 649 predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference 650 on knowledge discovery and data mining, pp. 1135–1144, 2016. 651 Robert A Rigby and D Mikis Stasinopoulos. Generalized additive models for location, scale and 652 shape. Journal of the Royal Statistical Society Series C: Applied Statistics, 54(3):507–554, 2005. 653 654 David Rügamer. mixdistreg: An r package for fitting mixture of experts distributional regression 655 with adaptive first-order methods. arXiv preprint arXiv:2302.02043, 2023. 656 657 David Rügamer, Florian Pfisterer, Bernd Bischl, and Bettina Grün. Mixture of experts distributional regression: implementation using robust estimation with adaptive first-order methods. AStA Ad-658 vances in Statistical Analysis, 108(2):351-373, 2024. 659 660 Mohammed Saeed, Mauricio Villarroel, Andrew T Reisner, Gari Clifford, Li-Wei Lehman, George 661 Moody, Thomas Heldt, Tin H Kyaw, Benjamin Moody, and Roger G Mark. Multiparameter 662 intelligent monitoring in intensive care ii (mimic-ii): a public-access intensive care unit database. 663 *Critical care medicine*, 39(5):952, 2011. 664 Sunil K Sapra. Generalized additive models in business and economics. International Journal of 665 Advanced Statistics and Probability, 1(3):64–81, 2013. 666 667 M Scott, Lee Su-In, et al. A unified approach to interpreting model predictions. Advances in neural 668 information processing systems, 30:4765–4774, 2017. 669 670 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, 671 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 672 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. 673 674 Sheng Shen, Zhewei Yao, Chunyuan Li, Trevor Darrell, Kurt Keutzer, and Yuxiong He. Scaling 675 vision-language models with sparse mixture of experts. In Findings of the Association for Com-676 putational Linguistics: EMNLP 2023, pp. 11329–11344, 2023. 677 678 Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 679 10781-10790, 2020. 680 681 Sarah Tan, Rich Caruana, Giles Hooker, Paul Koch, and Albert Gordo. Learning global additive 682 explanations for neural nets using model distillation. stat, 1050:3, 2018. 683 684 Anton Frederik Thielmann, René-Marcel Kruse, Thomas Kneib, and Benjamin Säfken. Neural 685 additive models for location scale and shape: A framework for interpretable neural regression beyond the mean. In International Conference on Artificial Intelligence and Statistics, pp. 1783-686 1791. PMLR, 2024. 687 688 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 689 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa-690 tion processing systems, 30, 2017. 691 692 Weibin Wu, Yuxin Su, Xixian Chen, Shenglin Zhao, Irwin King, Michael R Lyu, and Yu-Wing Tai. Towards global explanations of convolutional neural networks with concept attribution. In Pro-693 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8652-694 8661, 2020. 696 Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. 697 On completeness-aware concept-based explanations in deep neural networks. Advances in neural 698 information processing systems, 33:20554–20565, 2020. 699 Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learn-700 ing based natural language processing. *ieee Computational intelligenCe magazine*, 13(3):55–75, 701 2018.

702	Wei Zhang Brian Barr and John Paisley, Gaussian process neural additive models. Proceedings
703	of the AAAI Conference on Artificial Intelligence, 38(15):16865–16872, Mar. 2024. doi: 10.
704	1609/aaai.v38i15.29628.URL https://ojs.aaai.org/index.php/AAAI/article/
705	view/29628.
706	
707	
708	
709	
710	
711	
712	
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
720	
729	
720	
729	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

A DATASET DESCRIPTION

Housing Dataset. The California Housing dataset, cited from Pace & Barry (1997), encompasses a regression task based on median housing prices across California's census blocks. It comprises 20,640 instances, each characterized by 8 distinct attributes.

762 MIMIC-II Dataset. The MIMIC-II (Multiparameter Intelligent Monitoring in Intensive Care)
 763 dataset, as described by Saeed et al. (2011), facilitates a binary classification task aimed at pre 764 dicting mortality in intensive care units (ICUs). It contains 17 attributes, including 7 categorical
 765 variables.

MIMIC-III Dataset. MIMIC-III is a more extensive and detailed iteration of the MIMIC database
 referenced from Johnson et al. (2016). In our study, the same setting in NODE-GAM (Chang et al., 2022) is adopted, wherein categorical variables have been transformed into dummy variables for
 enhanced analysis.

771

766

Income Dataset. Originating from the UCI Machine Learning Repository (Blake, 1998), the Income dataset underpins a binary classification task. Its objective is to predict individuals earning in excess of \$50,000 annually.

Credit Dataset. The Credit dataset¹ provides samples for a binary classification task on transaction fraud detection. It contains 30 anonymized features including 28 coefficients of PCA components. In Credit, 492 out of 284,807 transactions are labeled as frauds.

Year Dataset. The Year dataset² contains data for a regression task, which uses the audio features to predict the release year of a song. It includes 515,345 samples with 90 features.

The statistics of each dataset can be found in Table 4.

783 784

779

781

Table 4: Statistics of real-world datasets used in model evaluation.

Dataset	Task	#Instances	#Features	#Categorical Features	#Classes
Housing	regression	20,640	8	0	_
MIMIC-II	classification	24,508	17	7	2
MIMIC-III	classification	27,348	57	0	2
Income	classification	32,561	14	8	2
Credit	classification	284,807	30	0	2
Year	regression	515,345	90	0	_

795 796

797

798

799

800

801

B IMPLEMENTATION DETAILS

For a comprehensive and equitable evaluation, we employed the officially released codes for baseline models including NODE, NODE-GAM/NODE-GA²M, EBM/EB²M, and NBM/NB²M. We adopted the PyTorch implementation of NAM/NA²M by NBM, which has been benchmarked against their model, for a direct comparison with our proposed models. In developing MixNAM, we constructed a multi-channel, fully connected network capable of encoding multiple features independently in a parallel manner.

For a fair comparison with existing methods, our experiments employ the processed version of MIMIC-II, MIIMC-II, Income, and Credit from (Chang et al., 2022), where the datasets are split into five parts for a five-fold cross-validation. The Housing dataset and the Year dataset are split into the training, validation, and test sets following setups of previous research (Chang et al., 2022;

807

^{807 &}lt;sup>1</sup>The Credit dataset can be downloaded at https://www.kaggle.com/datasets/mlg-ulb/ creditcardfraud.

²The Year dataset can be downloaded at https://archive.ics.uci.edu/dataset/203/ yearpredictionmsd

Radenovic et al., 2022), and each of them is tested with 10 different seeds to examine the variation of predictions. Consistent with previous studies (Chang et al., 2022; Popov et al., 2020; Radenovic et al., 2022), we applied the same quantile transformation to the features. Our evaluation metrics
include the Root Mean Square Error (RMSE) for the Housing dataset, the Area Under the Receiver
Operating Characteristic Curve (AUC) for the Income, MIMIC-II, and MIMIC-III datasets, and the
Mean Square Error (MSE) for the Year dataset following existing research (Chang et al., 2022;
Popov et al., 2020; Radenovic et al., 2022).

To be consistent with existing NAMs (Agarwal et al., 2021; Radenovic et al., 2022), the feature encoders f_1, \dots, f_n in MixNAM are implemented with MLPs. The expert predictors E_{11}, \dots, E_{nC} are implemented as one linear layer in our experiments as mentioned in Section 3.1. During training, we decrease the learning rate with cosine annealing following NBM (Radenovic et al., 2022). AdamW optimizer (Loshchilov & Hutter, 2018) is used to optimize the training objective.

- All experiments were run on NVIDIA A100 GPUs (40 GB / 80 GB).
- C HYPERPARAMETERS

822

823 824

825 826

827 828

829

830 831

832

833

834

835

836 837

838

839

840

841

842

843 844

845

846

847

848

849

850

851 852

853 854 855

856

We randomly search hyperparameters in the following ranges:

- #layers: the number of layers in each feature encoder, sampled from {3, 4}.
- Hidden dimension: the number of nodes in each layer of the feature encoder, sampled from {64, 128}.
 - #total experts: the number of total experts for each feature, sampled from {4}.
- #activated experts: the number of activated experts for each feature, sampled from {4}.
- Batch size: the number of samples in one batch during training, sampled from {512, 1024, 2048}.
- Max iteration: the number of iterations for training, sampled from {75, 150, 500, 1000}.
- Learning rate: the speed of gradient descent, sampled from [1e-6, 1e-1].
 - Weight decay: the coefficient for the L2 normalization on parameters, sampled from [1e-8, 1e-1].
 - Dropout: the probability of a parameter in feature encoders being replaced as 0 during training, sampled from {0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9}.
- Dropout expert: the probability of an expert's output being replaced as 0 during training, sampled from {0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9}.
- Output penalty: the coefficient for the L2 normalization on the outputs of features for reducing unimportant ones, sampled from [1e-8, 1e-1].
- Variation penalty: the weight for the expert variation penalty, sampled from $\{0.1\}$.
- Normalization: normalization methods used in the network, sampled from {batch_norm, layer_norm}.

The best hyperparameters we found for MixNAM are shown in Tables 5.

D MORE RELATED WORK ON GENERALIZED ADDITIVE MODELS WITH MIXTURE OF EXPERTS

In addition to the general related work on the design of generalized additive models (GAMs), plenty
of efforts have been made in the application of GAMs (Hastie & Tibshirani, 1995; Sapra, 2013;
Pedersen et al., 2019; Izadi, 2021), the evaluation of GAMs (Hegselmann et al., 2020; Chang et al.,
2021), and the generalization of GAMs to high-order feature interaction modeling (Enouen & Liu,
2022; Tan et al., 2018). In the context of Mixture of Experts (MoE), there is some specific related
work that tried to combine the idea of MoE with GAMs. Below we will introduce these attempts and discuss the differences between them and our MixNAM.

865		Table 5: H	yperparamete	ers for MixNA	AM on all dat	asets.	
866	Hyperparameter	Housing	MIMIC-II	MIMIC-III	Income	Credit	Year
867	#layers	4	4	4	4	4	4
868	Hidden dimension	128	128	128	128	128	128
869	#total experts	4	4	4	4	4	4
870	#activated experts	4	4	4	4	4	4
871	Batch size	2048	2048	1024	2048	2048	512
070	Max iteration	1000	1000	500	1000	150	75
872	Learning rate	5.97e-4	1.97e-4	1.68e-4	1.11e-4	4.00e-5	1.17e-4
873	Weight decay	5.29e-5	8.06e-4	2.78e-4	4.25e-3	1.88e-3	2.28e-6
874	Dropout	0.1	0.0	0.1	0.0	0.3	0.1
875	Dropout expert	0.2	0.4	0.2	0.5	0.2	0.6
876	Output Penalty	1.97e-5	3.49e-8	4.86e-5	3.99e-1	5.07e-4	1.45e-4
077	Variation Penalty	0.1	0.1	0.1	0.1	0.1	0.1
878	Normalization	layer_norm	layer_norm	layer_norm	layer_norm	batch_norm	layer_norm

GAMLSS (Rigby & Stasinopoulos, 2005) uses additive models to learn the distribution parameters of the target output such as mean and variance. While it is more flexible than standard GAMs, GAMLSS is still a strictly additive model. For example, consider a Gaussian output $y \sim \mathcal{N}(\mu, \sigma^2)$ modeled by $\mu = \sum_{i=1}^{n} f_i(x_i)$ and $\sigma = \sum_{i=1}^{n} g_i(x_i)$. With the reparameterization trick, the prediction can be re-written as

$$y = \mu + \sigma \varepsilon = \sum_{i=1}^{n} \left(f_i(x_i) + \varepsilon g_i(x_i) \right), \tag{15}$$

where x_1, \dots, x_n are the *n* features and $\varepsilon \sim \mathcal{N}(0, 1)$. It can be seen that the output is still a strict addition of different feature information. Similarly, FlexMix (Leisch, 2004) provides a mixture of GAMs where the expert weights are learned and fixed by an EM algorithm. These two models are still strictly additive in general.

Mixdistreg (Rügamer, 2023; Rügamer et al., 2024) models flexible mixture weights dynamically determined by the input. Actually, it can be considered as a simplified version of our MixNAM. In Mixdistreg, the mixture of GAMs can be modeled as

 $F(x) = \sum_{k=1}^{K} \left(\pi_k(x) \sum_{i=1}^{n} f_{ki}(x_i) \right) = \sum_{i=1}^{n} \left(\sum_{k=1}^{K} \pi_k(x) f_{ki}(x_i) \right),$ (16)

901 902 903

864

879 880

886 887

889 890 891

892

893

894

895

896

897

898 899

900

where K is the number of experts and $\pi_1(x), \dots, \pi_K(x)$ are the learned weights for each expert given x. While the weights are dynamically learned, they are kept the same across different features, which is less flexible than our MixNAM, where the learned weights could be diverse for different features. Also, the control of accuracy-interpretability trade-offs using our proposed variation penalty is not introduced in the previous work.

In general, classic methods such as GAMLSS (Rigby & Stasinopoulos, 2005) and FlexMix (Leisch, 2004) are strict additive models, which cannot mitigate the additivity constraint to achieve a balance between model performance and interpretability. Meanwhile, the recent Mixdistreg (Rügamer, 2023; Rügamer et al., 2024) mitigates the constraint but can be considered a simplified version of our MixNAM, with no control on the trade-offs between model accuracy and interoperability.

In addition to additive models, another prevalent approach for interpretability in tabular data involves
post-hoc feature attribution. Methods like LIME (Ribeiro et al., 2016) and SHAP (Scott et al., 2017) offer local explanations by detailing feature contributions for individual predictions. However,
MixNAM stands apart by providing a transparent, global understanding of how features influence predictions across the entire dataset.

918 E ADDITIONAL SIMULATION STUDIES ON EXTREME CASES

E.1 SIMULATION OF FEATURES WITH HIGH SPARSITY

In addition to the simulation study presented in Section 4.4, we further explore how MixNAM performs when processing data with highly sparse features. Following the original simulation settings, the output of the simulated multimodal data is computed as:

$$y = x_2 \sin(4\pi x_1) + x_2 + \varepsilon$$
 with $\varepsilon \sim \mathcal{N}(0, 0.1^2)$. (17)

Different from the previous simulation where x_2 is sampled from $\{-1, 1\}$ with equal probabilities, we manually control the proportion of data points with $x_2 = 1$ to represent 25%, 5%, and 1% of the dataset.

Figure 6 shows the learned shape plots of NAM and MixNAM on simulated data with different sparsities. We also explore how the number of experts C will affect the shape function learning of MixNAM. For simplicity, all experts will be activated during the expert routing (i.e., K = C). The results demonstrate that NAM tends to focus on the majority group when the signal for $x_2 = 1$ is sparse. In contrast, MixNAM successfully identifies both modalities in the data distribution and learns the sparse distribution better with an increasing number of experts.



Figure 6: Shape plots of MixNAM and NAM on the simulated data with high sparsity.

E.2 SIMULATION OF FEATURES WITH HIGH VARIABILITY

As we introduced above, the variability in model prediction refers to the extent to which the predicted outputs of a model fluctuate or differ when the same input feature value is presented across different samples. Formally, consider the mapping from the input features x_1, \dots, x_n to the predicted output:

$$\hat{y} = F(x_1, x_2, \cdots, x_n), \tag{18}$$

where F represents the underlying predictive function. For a fixed value of $x_i = a$, the variability of the contribution of x_i to \hat{y} can be defined as:

variability
$$|_{x_i=a} = \operatorname{Var}_{x_1, \dots, x_n} [F(x_1, x_2, \dots, x_n | x_i = a) - \operatorname{E}_{x_i} (F(x_1, x_2, \dots, x_n))].$$
 (19)

This definition measures how the contributions of $x_i = a$ deviate due to interactions with other features, encapsulating variability caused by feature dependencies. For traditional additive models, the variability defined above reduces to zero because the contributions of each feature are independent and do not interact. In models where interactions exist, this term captures the extent to which other features $x_j (j \neq i)$ influence the contribution of x_i .

To assess the performance of MixNAM under conditions of extreme variability in output contributions caused by a large number of modalities, we simulate a data distribution using the following equation:

$$y = \varepsilon + \sum_{i=2}^{NC+1} \frac{T}{NC} x_i \sin(4\pi x_1), \qquad (20)$$

where the feature x_1 is sampled from U(0, 1) and x_2, \dots, x_{NC+1} are the NC categorical features sampled from $\{-1, +1\}$ with equal probability. T is a scaling factor that amplifies variability across different modalities. For this study, we set T = 64 and evaluated both NAM and MixNAM across scenarios with NC = 1, 2, 4, 8, 16. The number of simulated samples is set dynamically as $2000 \times$ NC.

As illustrated in Figure 7, the results demonstrate that MixNAM effectively captures the multimodal contributions of x_1 to the output, showing consistent performance as the number of features and modalities increases. In contrast, NAM struggles to account for multimodality due to its inherent limitation of feature additivity. The results highlight the robustness and scalability of MixNAM in handling multimodal data with high variability.



Figure 7: Shape plots of MixNAM and NAM on the simulated data with high variability.

F ADDITIONAL ANALYSIS OF EXPERT ROUTING IN MIXNAM

While the expert routing mechanism in MixNAM can capture certain feature interactions and provide a continuous output range between the lower and upper bounds for each feature, it should be noted that the router is not a universal function approximator that dominates the model prediction even if there is no variation penalty.

Specifically, if each feature has only two expert functions to model the upper and lower bounds,
 both of which are activated, the relevant score estimation of each expert is a normalized Generalized
 Additive Model (GAM) with the sigmoid function as the link function. To prove this, we can rewrite
 the relevance score for the first expert of the *i*-th feature (Formula 6) as

1026

1027 1028

1030 1031

$$= \frac{1}{\exp(e_1^\top \varphi_i) + \exp(e_2^\top \varphi_i)} \\ \exp((e_1 - e_2)^\top \varphi_i)$$

 $r_{i1} = \frac{\exp(\varphi_i[1])}{\exp(\varphi_i[1]) + \exp(\varphi_i[2])}$

 $= \frac{\exp((e_1 - e_2) - \varphi_i)}{\exp((e_1 - e_2)^\top \varphi_i) + 1}$

 $\exp(e_1^{\top}\varphi_i)$

1032 1033 1034

1035 1036

1039 1040 1041

1043

1045 1046 1047

$$= \sigma((e_1 - e_2)^{\top}$$

$$=\sigma((e_1-e_2)^{\top}\varphi_i)$$

$$= \sigma \left((e_1 - e_2)^\top (\mu_i + \sum_{j=1}^n A_{ji}^\top f_j(x_j)) \right)$$
 (Formula 4)

$$= \sigma \left(e_1^{\top} \mu_i - e_2^{\top} \mu_i + \sum_{j=1}^n \left(e_1^{\top} A_{ji}^{\top} f_j(x_j) - e_2^{\top} A_{ji}^{\top} f_j(x_j) \right) \right)$$

With 1042

$$g_{i0} := e_1^\top \mu_i - e_2^\top \mu_i \quad \text{and} \quad g_{ij}(\cdot) := e_1^\top A_{ji}^\top f_j(\cdot) - e_2^\top A_{ji}^\top f_j(\cdot),$$

 $(e_1 := [1, 0]^{\top}, e_2 := [0, 1]^{\top})$

 $(\sigma(\cdot) := \text{Sigmoid}(\cdot))$

(21)

the relevance score can be further rewritten as 1044

$$r_{i1} = \sigma \Big(g_{i0} + \sum_{j=1}^{n} g_{ij}(x_j) \Big),$$

1048 which is a normalized GAM where the linked function is a sigmoid function. Thus, instead of interpolating all values between the upper and lower bounds, the expert routing mechanism can only 1049 interpolate certain positions whose values are restricted by a normalized GAM. 1050

1051 Similarly, for more general cases with K activated experts among C experts, the estimated relevance 1052 for the activated experts can be considered K probability values normalized by the softmax function. 1053 By definition in Formula 4, each probability before the normalization is given by a GAM.

1054 Such an analysis can be verified by the result of our additivity evaluation in Section 4.3.2, where the 1055 additivity of MixNAM is significantly higher than MLP, even when there is no variation penalty. 1056

1057 1058

1059

MIXNAM-D: MIXNAM WITH DIAGONAL ROUTING MATRIX G

As mentioned in Section 3.2, we explored a special version of MixNAM, termed MixNAM-D, where the scoring matrics compose a block diagonal matrix. In MixNAM-D, the matrix A_{ii} in Formula 1061 4 becomes a zero matrix if $i \neq j$. This configuration allows MixNAM-D to function similarly to 1062 NAMs, where the relevance estimation of a feature's experts is solely based on its own encoded 1063 information. 1064

Since the relevance for experts of each feature is solely determined by itself in MixNAM-D, the 1065 estimated relevance scores will be fixed for the same feature value. To encourage the model to learn various distributions with multiple experts, we introduce randomness into the decision-making 1067 process by employing the Gumbel-softmax technique (Jang et al., 2016) with the temperature $\tau =$ 1068 0.1 to re-sample the experts during training based on their estimated relevance scores in Formula 6: 1069

1070 1071

1072

$$\hat{r}_{ik} = \exp\left(\frac{\log(r_{ik}) + g_i}{\tau}\right) / \sum_{l=1}^{C} \exp\left(\frac{\log(r_{il}) + g_l}{\tau}\right), \quad \text{where } g_i, g_l \sim Gumbel(0, 1).$$
(22)

1074

The performance comparison of MixNAM-D with other baselines is presented in Table 6. The 1075 results demonstrate that, by regularizing the scoring matrix to be diagonal, MixNAM-D shows a 1076 performance close to traditional additive models, which have worse accuracies than complex models 1077 that capture feature interactions. 1078

Figure 8 shows the shape plots generated by MixNAM-D compared to other baseline additive mod-1079 els. As discussed above, MixNAM-D is a strictly additive model where the effect of each feature on

		Housing	MIMIC-II	MIMIC-III	Income	Credit	Year
Mode	el FA	$RMSE\downarrow$	AUC ↑	AUC ↑	AUC ↑	$\overline{\text{AUC}\uparrow}$	MSE ↓
MID	Y	0.501	0.835	0.815	0.914	0.981	78.48
WILI	<u>^</u>	± 0.006	± 0.014	± 0.009	± 0.003	± 0.007	± 0.56
VCD	aast V	0.443	0.844	0.819	0.928	0.978	78.53
AUD	Jost 🔨	± 0.000	± 0.012	± 0.004	± 0.003	± 0.009	± 0.09
FD ²	. v	0.492	0.848	0.821	0.928	0.982	83.16
EB N	1 ^	± 0.000	± 0.012	± 0.004	± 0.003	± 0.006	± 0.01
N1 A 2 N	π Υ	0.492	0.843	0.825	0.912	0.985	79.80
NA I	/I ^	± 0.008	± 0.012	± 0.006	± 0.003	± 0.007	± 0.05
EDM		0.559	0.835	0.809	0.927	0.974	85.81
EDM		± 0.000	± 0.011	± 0.004	± 0.003	± 0.009	± 0.11
NAM		0.572	0.834	0.813	0.910	0.977	85.25
INAN	•	± 0.005	± 0.013	± 0.003	± 0.003	± 0.015	± 0.01
MixN		0.553	0.830	0.805	0.927	0.977	85.60
MIXIN		± 0.001	± 0.012	± 0.006	± 0.003	± 0.010	± 0.04

Table 6: Comparison of MixNAM-D to other baselines on benchmark datasets. "FA" (Feature

the final output is determined by the feature itself. Thus, we plot MixNAM-D in the same way as
we did for EBM and NAM, directly showing the estimated contribution given by each feature value.
Moreover, since MixNAM-D is trained under the general MixNAM framework with multiple experts learned for each feature, we plot another figure for MixNAM-D including the estimated upper and lower bounds based on the learned experts.

It can be observed from Figure 8 that the patterns captured by MixNAM-D are similar to the ones by
EBM and NAM. Beyond the point estimation provided by traditional additive models, MixNAM-D
offers additional insights with its prediction bounds. In areas with low data density, the plot displays
a wide gap between the upper and lower bounds, highlighting variability in predictions where less
data is available.



Figure 8: Visual comparison of the "Longitude" feature impact on house prices by MixNAM-D.

H MIXNAM-E: MIXNAM WITH EVENLY DISTRIBUTED EXPERT ACTIVATION

While the activation strategy detailed in Section 3.2 facilitates continuous expert relevance estimation, we propose an adaptation that shifts this estimation to a discrete framework. This modification involves adjusting the weights to be evenly distributed across relevant experts, leading to a modified model termed MixNAM-E. To implement this, we use the masking vector outlined in Formula 5 and modify the relevance computation in Formula 6 as follows:

$$r_{ik} = \exp(\varphi_i[k] - \varphi_i^*[k] + M_i[k]) / \sum_{l=1}^C \exp(\varphi_i[l] - \varphi_i^*[l] + M_i[l]).$$
(23)

 φ_i^* mirrors the values of φ_i but does not require gradient computations. The activation strategy implemented by Formula 23 will result in a finite set of possible outcomes for each input feature value. Given a configuration of C experts and K activated ones per feature, the predicted output by MixNAM-E for any given feature will be one of the possible selections from $\binom{U}{K}$ combinations.

Table 7 shows the performance (RMSE) of MixNAM-E on the Housing dataset with different varia-tion penalties, along with their corresponding shape plots of the "Longitude" feature. All the tested MixNAM-E models have a configuration of C = 32 and K = 16. In the plots, we illustrate the upper and lower bounds for each feature by selecting the maximum and minimum values from the $\binom{C}{K}$ possible combinations.

Compared with MixNAM (Table 2), the performance of MixNAM-E is slightly worse due to the limitations imposed by its discrete range. However, even without any variation penalty, the esti-mated bounds by MixNAM-E still fit the actual score distribution tightly, accurately reflecting the prediction limits.



Table 7: Performance (RMSE \downarrow) and explanations of MixNAM-E with different variation penalties.



PERFORMANCE OF MIXNAM WITH THE SCALING OF EXPERTS Ι

As MixNAM employs a mixture of experts for data modeling, we explore how the number of total experts (C) and the number of activated experts (K) affect the overall model performance. Table 8 presents the performance of both MixNAM and MixNAM-E on the Housing dataset under vari-ous configurations. The results indicate that MixNAM reaches optimal performance at C = 8 and K = 4, and further increases in C or K do not enhance its performance. Such a finding is consistent with the role of experts in MixNAM, which primarily explore the variability in the prediction space to learn the upper and lower bounds, as the dynamic routing mechanism offers a continuous range within the bounds for prediction. In contrast, MixNAM-E benefits from an increased number of C and K, showing improved performance as these parameters grow. This difference shows how the continuous expert activation in MixNAM contrasts with the discrete one in MixNAM-E, highlighting the flexibility of our MixNAM framework to adapt to various scenarios through different implementations.

Table 8: Results for different numbers of activated experts (K) and all experts (C).

					1				
		Mixl	NAM		MixNAM-E				
$K \setminus C$	2	4	8	16	16	32	64	128	
2	0.479	0.459	0.455	0.460	0.486	0.487	0.486	0.489	
2	± 0.003	± 0.003	± 0.004	± 0.004	± 0.007	± 0.005	± 0.006	± 0.005	
4		<u>0.451</u>	<u>0.447</u>	0.455	0.469	0.466	0.468	0.470	
4	_	± 0.002	± 0.004	± 0.004	± 0.007	± 0.005	± 0.004	± 0.004	
8			<u>0.449</u>	0.456	0.470	0.459	<u>0.458</u>	0.458	
0	_	_	± 0.004	± 0.003	± 0.003	± 0.004	± 0.002	± 0.004	
16				0.456	0.587	0.462	<u>0.451</u>	<u>0.453</u>	
10	_	_	_	± 0.005	± 0.003	± 0.002	± 0.002	± 0.002	

1188 J DISCUSSION ON MODEL COMPLEXITY

1189 1190

As MixNAM is a fundamental framework aiming to extend traditional additive models, our design 1191 prioritizes the performance and comprehensiveness of the overall framework over the optimization 1192 of model efficiency. On the basis of Neural Additive Models (NAMs), MixNAM introduces more 1193 parameters in the expert encoding and dynamic routing steps. To rigorously analyze how much 1194 cost MixNAM will bring compared to NAM, suppose we have n features, d-dimensional output for feature encoding, and C total experts for each feature. In our experiments, MixNAM is implemented 1195 with d = 128 and C = 4. For MixNAM-D we have d = 128 with C = 64. The theoretical and 1196 empirical additional costs brought by MixNAM and its variant are presented in Table 9, including 1197 the increase in parameter count and the additional memory required (assuming float32 storage). 1198

1199

1201

Table 9: Theoretical and empirical additional costs brought by MixNAM compared to NAM. Mix-NAM is implemented with d = 128, C = 4. MixNAM-D is implemented with d = 128, C = 64.

	#Features	MixN	MixNAM-D		
		Parameter Count	Memory Usage	Parameter Count	Memory Usage
Housing	8	37k	144k	132k	516k
MIMIC-II	17	157k	613k	281k	1.1M
MIMIC-III	57	1.7M	6.5M	941k	3.59M
Income	14	108k	420k	231k	903k
Credit	30	476k	1.8M	495k	1.89M
Year	90	4.2M	16.0M	1.5M	5.67M
Theoretical	-	nC[(n+1)]	1)d + 2]	nC(2d)	+2)

1212

The results indicate that the additional cost associated with MixNAM-D scales linearly with the 1213 number of features, whereas the cost for MixNAM includes a squared term due to its more complex 1214 routing mechanism. Specifically, MixNAM utilizes a full $n \times n$ block matrix to encode feature 1215 interactions, which is simplified as a block diagonal matrix in MixNAM-D. Despite this increased 1216 complexity, the additional memory usage remains manageable with current computing resources. 1217 Future efforts could be made to sparsify the routing matrix in MixNAM, which could potentially 1218 reduce the additional cost while retaining the model performance.

1219 1220 1221

1222

MORE VISUALIZATION RESULTS OF MIXNAM ON REAL-WORLD DATA Κ

Figures 9 - 12 show the complete visualization results of feature explanations by MixNAM on 1223 different datasets. Here we present results for datasets with no more than 50 features, including 1224 Housing, MIMIC-II, Income, and Credit. 1225

1226 1227

1228

DISCUSSIONS, LIMITATIONS AND FUTURE WORK L

1229 MixNAM can have a significant impact on high-stakes areas such as finance and healthcare where 1230 the explanation of predictions is as crucial as the predictions themselves. The introduction of Mix-1231 NAM advances the capabilities of additive models, which may encourage the deployment of re-1232 lated research in practical applications. Moreover, MixNAM generalizes existing Neural Additive 1233 Models, offering a new direction to enhance additive models by balancing the accuracy and interpretability under the MixNAM framework. However, MixNAM does not inherently address issues 1234 of fairness or mitigate implicit bias in the data. Instead, it provides transparent feature explanations 1235 that may highlight these biases. Therefore, careful considerations and responsible use are essential 1236 when applying MixNAM in scenarios where fairness is a critical concern. The transparency of Mix-1237 NAM may also serve as a tool for identifying data biases, contributing to equitable and accountable AI systems. 1239

We acknowledge that this work has certain limitations, which could be the direction of future re-1240 search. First, while MixNAM provides a comprehensive visualization of feature impacts on predic-1241 tions, the intervenability of additive models is reduced due to the uncertainty in outcome predictions



Figure 9: Visualization of feature explanations by MixNAM on the Housing dataset.

when altering single feature values. However, MixNAM compensates by providing a range of pos-sible outcomes, defined by estimated upper and lower bounds. The expert variation penalty can also adjust trade-offs between intervenability and performance. Future efforts could be made to handle the uncertainty in MixNAM to offer precise estimations after interventions. Second, although recent advancements have improved the efficacy and scalability of Neural Additive Models (Radenovic et al., 2022; Zhang et al., 2024), our MixNAM is constructed based on the foundational vanilla NAM, which serves as the most general framework for neural-network-based additive models. En-hancing the efficiency of MixNAM could be a valuable direction for subsequent research. Third, we have tested our MixNAM only on tabular data, since features are clearly defined in this modality, aligning with prior studies on additive models (Agarwal et al., 2021; Chang et al., 2022; Lou et al., 2012). However, our framework has the potential to be generalized to other modalities. For example, features in images can be extracted with object detection (Bochkovskiy et al., 2020; Tan et al., 2020) or concept learning methods (Ghorbani et al., 2019), and the feature encoders in MixNAM could be implemented with CNN (He et al., 2016; Krizhevsky et al., 2012) or Transformer (Dosovitskiy et al., 2020; Vaswani et al., 2017) to process data in different modalities. We leave this as one potential direction for future research and applications.









Figure 12: Visualization of feature explanations by MixNAM on the Credit dataset.