# A GENERALIZED ADDITIVE MODELS: EXTENDED RELATED WORK

As with many families of machine learning algorithms, the differences among GAM algorithms lie in (a) the functional form of the shape functions  $f_i$ , (b) the learning algorithm used for their estimation and (c) regularity assumptions and regularization. Two important properties that all GAMs share are (1) the ability to learn non-linear transformations for each feature and (2) additively combining these shape functions (prior to applying the link function) to create modularity that aids interpretability by allowing users to examine shape functions one-at-a-time.

710 Typically, GAMs have relied on splines and backfitting algorithms for estimation (Hastie and Tibshi-711 rani, [1987], with subsequent works focusing on improving efficiency and stability through penalized 712 regression splines (Wood, 2003) and fast, stable fitting algorithms (Wood, 2001). Spline-based 713 GAMs are typically fitted using the backfitting algorithm, an iterative procedure that starts with initial estimates of the smooth functions for each predictor variable. The algorithm then repeatedly 714 updates each function by fitting a weighted additive model to the residuals of the other functions until 715 convergence is achieved. The weights are determined by the current estimates of the other functions 716 and the link function in the case of generalized additive models. 717

718 Modern approaches leverage machine learning advances. Explainable Boosting Machines (EBMs) 719 (Lou et al., 2012; 2013; Caruana et al., 2015) model the shape functions using decision trees, which 720 are fitted using a variant of gradient boosting called cyclic gradient boosting. The model iteratively learns the contribution of each feature and interaction term in a round-robin fashion, using a low 721 learning rate to ensure that the order of features does not affect the final model. This cyclic training 722 procedure helps mitigate the effects of colinearity among predictors by providing opportunity for 723 data-driven credit attribution among the features while preventing multiple counting of evidence. 724 EBMs are also popular because they can accurately capture steps in the shape functions, which is 725 important for modeling discontinuities in data, such as treatment effects in medical data. 726

More recently, Neural Additive Models (NAMs) (Agarwal et al.) [2021) and follow up works (Chang et al.) [2021; [Dubey et al.] [2022; [Radenovic et al.] [2022; Xu et al.] [2022; Enouen and Liu, [2022; Bouchiat et al.] [2024) use multilayer perceptrons (MLPs), as non-linear transformations, to model the shape functions  $f_i$ . As a result, NAMs can be optimized using variants of gradient descent by leveraging automatic differentiation frameworks.

Finally, GAMs have also found applications in time-series forecasting, with models such as
Prophet (Taylor and Letham, 2018) and NeuralProphet (Triebe et al., 2021). Interestingly, the
1-layer versions of the recently proposed Kolmogorov-Arnold Networks (KANs) (Liu et al., 2024)
may be viewed as GAMs with spline based shape functions.

736 737

738

## **B** DATASET DETAILS

In this section, we provide details on the datasets used in our empirical evaluations of GAMformerand other baselines in Section 4 of the main paper.

- 741 742
- B.1 TABPFN TEST DATASETS

As test dataset, we used the 30 datasets used in Hollmann et al. (2023) which were obtained from OpenML (Vanschoren et al., 2014). These were chosen because they contain up to 2000 samples, 100 features and 10 classes, show in Table 1.

746 747 748

752

B.2 BINARY CLASSIFICATION

749 Churn dataset. The Telco Customer Churn Dataset is a binary classification dataset for predicting
 750 potential subscription churners in a telecom company, containing customer information and churn 751 related features.

Adult dataset. The Adult dataset Dua and Graff (2017), also known as the "Census Income" dataset, is a widely-used benchmark for binary classification, predicting whether an individual's annual income exceeds \$50,000 based on 14 attributes from the 1994 United States Census Bureau data.

did	name	d	n	k	d	lid	name	d	n	k
11	balance-scale	5	625	3	104	49	pc4	38	1458	2
14	mfeat-fourier	77	2000	10	10	50	pc3	38	1563	2
15	breast-w	10	699	2	10	63	kc2	22	522	2
16	mfeat-karhunen	65	2000	10	10	68	pc1	22	1109	2
18	mfeat-morphological	7	2000	10	14	62	banknote-authentication	5	1372	2
22	mfeat-zernike	48	2000	10	14	64	blood-transfusion	5	748	2
23	cmc	10	1473	3	14	80	ilpd	11	583	2
29	credit-approval	16	690	2	14	94	qsar-biodeg	42	1055	2
31	credit-g	21	1000	2	15	10	wdbc	31	569	2
37	diabetes	9	768	2	63.	32	cylinder-bands	40	540	2
50	tic-tac-toe	10	958	2	233	81	dresses-sales	13	500	2
54	vehicle	19	846	4	409	66	MiceProtein	82	1080	8
188	eucalyptus	20	736	5	409	75	car	7	1728	4
458	analcatdata_authorship	71	841	4	409	82	steel-plates-fault	28	1941	7
469	analcatdata_dmft	5	797	6	409	94	climate-model	21	540	2

Table 1: Test dataset names and properties, taken from Hollmann et al. (2023). Here *did* is the OpenML Dataset ID, *d* the number of features, *n* the number of instances, and *k* the number of classes in each dataset.

Table 2: Comparison of GAMformer with other GAM variants and full complexity models on various datasets. We report ROC-AUC (%) (higher is better) and the standard error over 10 fold cross-validation. We also report results by pyGAM (Servén and Brummitt, 2018).

	GAMs	Full Complexity					
	GAMformer (ours)	EBM (Main effects)	Logistic Regression	pyGAM (Main effects)	EBM	XGBoost	Random Forest
Churn	$81.69 \pm 0.1$	$83.59 \pm 0.1$	81.66 ± 0.1	$82.03 \pm 0.0$	$83.68 \pm 0.1$	$83.53 \pm 0.0$	$82.07 \pm 0.0$
Support2	$80.84 \pm 0.1$	$82.36 \pm 0.0$	$81.1 \pm 0.0$	$81.74 \pm 0.2$	$83.51 \pm 0.0$	$84.03 \pm 0.0$	$83.93 \pm 0.0$
Adult	$90.05 \pm 0.0$	$93.05 \pm 0.0$	$90.73 \pm 0.0$	91.55 ± 0.0	$93.07 \pm 0.0$	$93.16 \pm 0.0$	$91.8 \pm 0.0$
MIMIC-2	$82.22 \pm 0.0$	$85.15 \pm 0.0$	$81.62 \pm 0.0$	$83.89 \pm 0.1$	$86.36 \pm 0.1$	$87.29 \pm 0.0$	$87.31 \pm 0.0$
MIMIC-3	$74.41 \pm 0.1$	$81.14 \pm 0.0$	$78.05 \pm 0.0$	79.95 ± 0.1	$82.52 \pm 0.1$	$83.32 \pm 0.0$	$81.28 \pm 0.1$

785 786 787

788

789

790

776

777

778

779

**MIMIC-II dataset.** The MIMIC-II dataset Lee et al. (2011b) is a publicly-available database of clinical data from diverse ICU patients, integrating demographics, vital signs, lab results, medications, procedures, notes, and imaging reports, along with mortality outcomes.

MIMIC-III dataset. The MIMIC-III dataset Johnson et al. (2016) expands on MIMIC-II, with a larger patient cohort, more recent records, enhanced data granularity, and the inclusion of free-text imaging report interpretations.

SUPPORT2 dataset. The SUPPORT2 dataset Connors Jr et al. (1996) contains medical information from critically ill hospitalized adults, compiled to study the relationships between medical decision-making, patient preferences, and treatment outcomes, with variables spanning demographics, physiology, diagnostics, treatments, and survival/quality of life outcomes.

799 800

801

# C PROPERTIES OF GAMFORMER

802 C.1 DATA SCALING

To assess GAMformer's ability to generalize to datasets containing more datapoints than it saw during training, i.e. larger context sizes, we conducted an experiment that varied the number of training data points and evaluated the impact on ROC-AUC performance using a consistent validation split. To ensure the robustness of our findings, we sampled training datasets three times with replacement for each training size. The results in Figure 8 demonstrate that GAMformer's ROC-AUC improves across datasets when the number of training examples is up to twice the number of training examples seen during training. For comparison, we also evaluated the performance of EBMs under the same



Figure 8: Demonstration of the ability of GAM former to scale beyond the datapoints seen during training while leveraging the additional data points to increase its performance. The dashed vertical line denotes the number of in-context examples seen during training (500).



Figure 9: Comparison of GAM former and EBMs in terms of (a) performance on class imbalanced data and (b) robustness to noisy labels. The shaded areas represent the 5% and 95% confidence intervals estimated using 1000 bootstrap samples.

833 834

837

839

840

818

819

820

821 822

823

824

825

827

828

829 830

831

832

835 conditions. While EBMs also exhibited improvements in ROC-AUC with increased training data, 836 they achieved higher accuracy when provided with a larger number of examples. This observation highlights a limitation of GAM former in its ability to fully leverage additional training samples. 838

C.2 **CLASS IMBALANCE** 

841 To compare GAM former's sensitivity to class imbalance with that of EBMs, we conduct the following analysis. First, we sample 300 data points from two centroids in a 20-dimensional feature space, 842 creating a binary classification problem. We then vary the ratio of the two classes to introduce 843 increasing levels of imbalance in the sampled data. Next, we split the data into train and test sets 844 using a 75% to 25% split and evaluate the performance using the AUC-ROC metric. We repeat the 845 experiment 10 times for each data ratio. Our results are shown in Figure 9a, the shaded area are the 846 5%, 95% confidence intervals estimated using 1000 bootstrap samples. We see that GAMformer 847 performs on average better than EBMs in this setting and shows no inherent sensitivity to class 848 imbalance. 849

C.3 NOISE ROBUSTNESS

852 To gain a deeper understanding of GAM formers' sensitivity to noisy or incorrect labels, we conducted 853 an experiment similar to the one described in Appendix C.2. We generated 300 data points and 854 randomly perturbed the labels in the train split with increasing probability (75%, 25% train/test split), 855 repeating each experiment 10 times. Figure 9b illustrates our findings. Once again, we observed that GAMformer exhibits a sensitivity to noisy labels comparable to that of EBMs. 856

857 858

859

850

851

#### D SYNTHETIC DATA PRIORS

We use the same synthetic data generation process proposed in Prior-Data-Fitted Networks 861 (PFNs) (Hollmann et al.) 2023; Müller et al., 2022) and provide a brief summary of the process. 862

TabPFN is trained on two synthetic data priors, which are mixed during training. TabPFN introduced 863 a synthetic data prior based on Structural Causal Models (SCMs). SCMs are particularly suitable for 864 modeling tabular data as they capture causal relationships between columns, a strong prior in human 865 reasoning. An SCM comprises a set of structural assignments (mechanisms) where each mechanism 866 is defined by a deterministic function and a noise variable, structured within a Directed Acyclic Graph 867 (DAG). The causal relationships are represented by directed edges from causes to effects, facilitating 868 the modeling of complex dependencies within the data. To instantiate a PFN prior based on SCMs, one defines a sampling procedure to create supervised learning tasks. Each dataset is generated from a randomly sampled SCM, including its DAG structure and deterministic functions. Nodes in the 870 causal graph are selected to represent features and targets, and samples are generated by propagating 871 noise variables through the graph. This process results in features and targets that are conditionally 872 dependent through the DAG structure, capturing both forward and backward causation (Hollmann 873 et al., 2023). This allows for the generation of diverse datasets. 874

The second prior samples of synthetic data using Gaussian Processes (GPs) (Rasmussen and Williams, 2006) with a constant mean function and a radial basis function (RBF) kernel to define the covariance structure. Hyperparameters such as noise level, output scale, and length scale are sampled from predefined distributions to introduce variability. Depending on the configuration, input data points can be sampled uniformly, normally, or as equidistant points and the target column is generated by passing the input data through the GP. This prior gives the model the ability to learn smoother functions.

For multi-class prediction, scalar labels are transformed into discrete class labels by partitioning the
 scalar values into intervals corresponding to different classes, ensuring the synthetic data is suitable
 for imbalanced multi-class classification tasks.

Finally, both priors are combined by sampling batches of data from each prior with different probabilities during training. In all of our experiments we sampled from the SCM and GP prior with
probability 0.96 and 0.04, respectively.

- 888
- 889 890

#### E TRAINING DETAILS

891 892

893

894

895

896

897

898

899

900

In GAMformer, we used a transformer model with 12 hidden layers, 512 embedding size and 4 heads per attention. To bin the shape functions and all features we used 64 bins. For training, we use the AdamW (Loshchilov and Hutter, 2019) optimizer ( $\beta_1 = 0.9$ ) and cosine learning rate schedule with initial learning rate of 3e-5, 20 warm up epochs and minimum learning rate of 1e-8 for 25 days on a A100 GPU with 80Gb of memory. We used mixed precision training. Each epoch (arbitrarily) consists of 65536 synthetic datasets; the model trained for 1800 epochs, meaning it saw over 100M synthetic datasets. We used a batch size of 8, that we doubled at epoch 20, 50, 200 and 1000. Each synthetic dataset consisted of 500 samples that were split into training and test portions using using a uniform sampling of the training fraction, and used a number of features drawn uniformly between 1 and 10.

905 906

907

908 909

910 911

## F HIGHER-ORDER EFFECTS

To handle higher-order effects, we compute the best pairs with the FAST algorithm (Lou et al., 2013) and evaluate GAM former on the top pairs using the following ratios of features:

 $\mathcal{P} = [0.01p, 0.05p, 0.1p, 0.2p, 0.4p, 0.8p, 0.9p]$ 

where we recall that p denotes the number of features. We round off each ratio to determine the number of target pair features, evaluate performance on hold-out validation data from the training set, and select the number of pairs with the best validation performance. The model is then fitted on the entire training dataset. This involves doing  $|\mathcal{P}| + 1$  forward passes, which is unproblematic as doing one forward pass is very fast, even on a CPU. One could also vectorialize all computations which we do not do given the low fitting time.

#### 918 G SHAPE FUNCTIONS

 In this section, we show complementary results on the shape functions estimates from GAMformer and EBM (main effects only) on the MIMIC-II (Lee et al.) 2011a) (complementary to the plots in Figure 7) and on the MIMIC-III datasets.

## G.1 MIMIC-II DATASET



Figure 10: The remaining shape functions derived from GAMformer and EBMs on the **MIMIC-II dataset** for critical clinical variables. The plot above each figure shows the data density. There are interesting differences between the EBM and GAMformer shape plots for several of the categorical variables. Although different GAM algorithms do not usually learn identical functions, we are investigating to better understand these differences.

#### G.2 MIMIC-III DATASET

 972 973 974 0.050 1.0).2GAMformer 0 0.8 Log-Odds (GAMFormer) 0.05 0.3 975 0.20.02 0.6 Sports (EBM) 0.4 Odds 976 0.0 0.0 0.00 0.0 0.0 0.0 0.000 977 0.0250.20.0 0.3 -0.1 0.2-0.5 -0.2978 0.050 0.0 1.0 adult'icu 0 1 admType'URGENT 0 1 admType'ELECTIVE admType<sup>\*</sup>EMERGENCY 979 980 0. 1.0 0.25).50 GAMformer 0.500.2 981 Log-Odds (GAMFormer) 0.4 0. EBM 0.20.00 0.250.250.50.1 Log-Odds (EBM) 982 0.2 0.0 -0.250.00 0.00 0.0 0.0 983 0.0 0.0 -0.25 -0.50 -0.25-0.3 0. 984 0.2 0.5 -0.75 -0.2-0.502040 985 2550 75 0.02.55.020 40 age albumin aniongap bicarbonate 986 .50 0.5987 GAMformer .0 0.20.20.25Log-Odds (GAMFormer) 0.2 EBM 0 0.2988 0.00.5spbO-go. 0.00 EBM) 0.1 0.0 989 0. 0.0 0.0 -0.25-0.5 0.0 -0.2 0.0 990 -0.50 0.4 -0.5 0. 0.1 991 -1.050 bilirubin 100 200 20 creatinine 100 40 chloride 992 bun 993 L.0 0.20.75GAMformer 0.4 0.5994 EBM Log-Odds (GAMFormer) 0. 0.50 0.50.1Log-Odds (EBM) 995 0.2 0.0 0.0 0.250.0 0.0 0.0 996 0 0.0 0.00 **Will** -0. -0.1 997 -0.1.0 5 -0.5 0.250.2 200 300 100 50 100 998 diashn'max diasbp'mean diasbp'min eth'asian 999 0.2 GAMformer 1000 0.20.050. 0.2Log-Odds (GAMFormer) 0.0  $_{\rm EBM}$ 0. 0.20.3 1001 Log-Odds (EBM) 0.0 0.0 0.0 0.0 0.0 0.0 0.00 0.00 1002 -0.1 1003 0 -0.20.05 -0.20. -0.05 -0.2 -0.2 1004 1 1 0 0 1 0 1 eth<sup>•</sup>black eth hispanic eth<sup>•</sup>other eth<sup>•</sup>white 1005 1006 GAMformer 0.0 0.05 Log-Odds (GAMFormer) 1007 0.05 0.2 0.20.4EBM Log-Odds EBM) 1008 0.20.0 0.00 0.000.00 0.0 0.0 ( 1009 0.0 -0.0 -0.2-0.05 -0.1 1010 0.0 -0.2 0.2 -0.41000 250 500 1011 first hosp'stay first'icu'stay glucose'max glucose'mean 1012 0.6 0. GAMformer 1013 0.2 0.25Log-Odds (GAMFormer) 0.40.2 EBM 0.50.20.00 0.5Log-Odds (EBM) 1014 0.5 0.10.2-0.25 1015 0.0 0.0 0.0 0.0 0.0 0.0 0.0 -0.50 1016 -0.2-0.5 -0. -0.75 -0.51017 -0.21000 2000 100 200 50 100 100 glucose heartrate max heartrate mean heartrate min 1018

Figure 11: The shape functions derived from GAMformer and EBMs on the **MIMIC-III dataset** for critical clinical variables. The plot above each figure shows the data density. The results are based on 30 models for both GAMformer and EBMs, each fitted on 10,000 randomly selected data points. There are interesting differences between the EBM and GAMformer shape plots for several of the categorical variables. Although different GAM algorithms do not usually learn identical functions, we are investigating to better understand these differences.



Under review as a conference paper at ICLR 2025

Figure 12: The remaining shape functions derived from GAMformer and EBMs applied to the
 MIMIC-III dataset for critical clinical variables. The plot above each figure shows the data density
 in the training set. The results are based on 30 models for both GAMformer and EBMs, each fitted
 on 10,000 randomly selected data points.

1079