

How Interpretable and Trustworthy are GAMs?

Chun-Hao Chang
University of Toronto, Vector
Institute, Hospital of Sick Children
kingsley@cs.toronto.edu

Sarah Tan
Cornell University
ht395@cornell.edu

Ben Lengerich
MIT, Broad Institute
blengeri@mit.edu

Anna Goldenberg
University of Toronto, Vector
Institute, Hospital of Sick Children
anna.goldenberg@utoronto.ca

Rich Caruana
Microsoft Research
rcaruana@microsoft.com

ABSTRACT

Generalized additive models (GAMs) have become a leading model class for interpretable machine learning. However, there are many algorithms for training GAMs, and these can learn different or even contradictory models, while being equally accurate. Which GAM should we trust? In this paper, we quantitatively and qualitatively investigate a variety of GAM algorithms on real and simulated datasets. We find that GAMs with high feature sparsity (only using a few variables to make predictions) can miss patterns in the data and be unfair to rare subpopulations. Our results suggest that inductive bias plays a crucial role in what interpretable models learn and that tree-based GAMs represent the best balance of sparsity, fidelity and accuracy and thus appear to be the most trustworthy GAM models.

CCS CONCEPTS

• Computing methodologies → Model verification and validation.

KEYWORDS

Generalized Additive Models, Interpretability, Inductive Bias

ACM Reference Format:

Chun-Hao Chang, Sarah Tan, Ben Lengerich, Anna Goldenberg, and Rich Caruana. 2021. How Interpretable and Trustworthy are GAMs?. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3447548.3467453>

1 INTRODUCTION

As the impact of machine learning on our daily lives continues to grow, we have begun to require that ML systems used for high-stakes decisions (e.g., in healthcare, finance and criminal justice) not only be accurate but also satisfy other properties such as fairness or interpretability [7, 18]. Generalized additive models (GAMs) have emerged as a leading model class that is designed to be accurate, and

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '21, August 14–18, 2021, Virtual Event, Singapore

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-8332-5/21/08...\$15.00
<https://doi.org/10.1145/3447548.3467453>

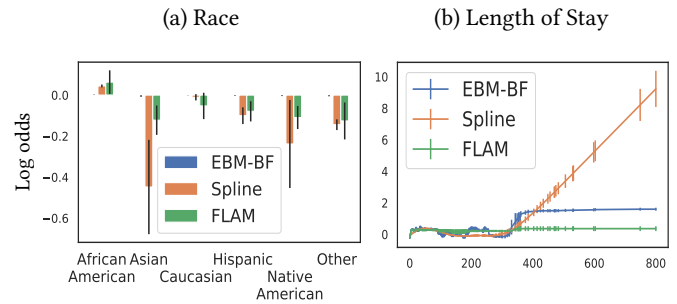


Figure 1: Three GAMs with similar accuracy trained on the COMPAS recidivism dataset (two of six features shown) that tell very different stories about the dataset. On the left, EBM-BF disagrees with FLAM and Spline about the presence of racial bias in the data. On the right, FLAM suggests that length of stay has no impact on risk, but Spline shows that risk grows strongly with length of stay.

yet simple enough for humans to understand and mentally simulate how a GAM model works [13], and is widely used in scientific data exploration [11, 16, 24] and model bias discovery [36, 37].

GAMs were originally trained using smoothing splines [10, 41] that enforced smoothness in the learned functions. Later, several trend-filtering based methods including fused lasso additive models were proposed to make learned functions more sparse and jumpy [29, 39]. Lou et al. [19] also proposed using boosted-tree-based methods to fit GAMs. Subsequent work showed the value of tree-based GAMs on two healthcare datasets [4], and also to help audit black-box models to ensure fairness [37].

Do GAMs trained with different algorithms agree with each other? In Fig. 1, we show that three GAMs with similar accuracy provide very different interpretations of the COMPAS recidivism dataset, a dataset in which bias is an important concern. For instance, in Fig. 2(a) EBM-BF suggests that there is no racial bias in the data, while Spline indicates that there is strong racial bias, yet is slightly less accurate. Should we believe EBM-BF because of its slightly higher accuracy and believe there is no racial bias? Probably not. Then how should we determine which GAM to believe?

In this paper, we try to answer this question by studying two aspects of GAMs trained with different algorithms. First, we quantify which GAMs use fewer features to make predictions (similar to

ℓ_1 -regularization), which we call *feature sparsity*. Although feature sparsity is sometimes preferred because it appears to yield simpler explanations [7, 38], it can be dangerous for data exploration as it can hide bias in the data. Consider a GAM that appears to be unbiased by showing no effect on sensitive variables such as race, but instead, because the learning algorithm is biased to use fewer features, it has compiled the racial bias into other correlated variables like zip code that are not obviously related to race, thus allowing the racial bias to go unrecognized. Furthermore, for features that only matter for rare subpopulations (e.g. a rare disease), a sparse-feature GAM could easily ignore such features but still remain accurate, leading to failure or discrimination for that subpopulation. In this paper, we empirically verify this phenomenon by showing that sparse GAMs often have higher loss for minority classes than less sparse GAMs.

Second, we examine how much we can trust each GAM to reflect true patterns in the data, a property we call *data fidelity*. Shmueli [33] contrasted predictive models that seek to minimize the *combination* of variance and bias (defined in a statistical sense, not in the sense of unfairness) to explanatory models that aim to capture true patterns in data by minimizing bias alone. For the former, bias can be sacrificed for improved variance, and Shmueli [33] provided examples of how the “wrong” model can sometimes predict better than the right one. In this paper, we study this phenomenon across different GAM algorithms. For real data where we do not know the underlying data patterns, we use the bias term from bias-variance analysis as a proxy for data fidelity. We also experiment with simulated datasets that have different data generators, each of which may favor GAM algorithms with certain inductive biases, and measure the worst-case data fidelity of each GAM algorithm across multiple datasets. This allows us to quantify if some GAM algorithms have high accuracy but low data fidelity, which may mislead users to trust the wrong explanations.

Our key contributions in this paper are:

- We compare different GAM algorithms on ten classification datasets and find that the most accurate GAMs yield similar accuracy, yet learn qualitatively different explanations.
- We measure which GAM algorithms lead to models that are more or less sparse, a property we call *feature sparsity*. We show that sparse-feature GAMs can discriminate on rare subpopulations leading to unfairness.
- We examine several case studies of data anomaly discovery to see which GAMs can or cannot be trusted to discover these true patterns in the data, a property we call *data fidelity*.
- We show that some GAMs have high accuracy but low data fidelity which can mislead users who select models by accuracy alone.
- We find that inductive bias plays a crucial role in model explanations, and recommend tree-based GAMs over other GAMs for their low feature sparsity and superior data fidelity.

2 RELATED WORK

While we study GAMs in this paper, GAMs are not the only interpretable model class to come under scrutiny recently. The instability of decision trees (another model class commonly considered interpretable) has been pointed out [9], and the vulnerability of post-hoc

explanation methods such as LIME [27] and Shapley values [20] to input perturbation has been exploited to generate adversarial attacks on model explanations [34]. Hooker and Mentch [14] also found partial dependence and feature importance metrics based on permuting inputs to be particularly misleading when inputs are highly dependent, and earlier work found feature importance metrics to be biased for certain types of models with different inductive biases, e.g., random forest feature importance is biased towards variables with many potential splits such as categorical variables with many levels [35, 45].

Our paper is not the first to compare different GAM algorithms, but to the best of our knowledge it is the first to focus on interpretability and its relationships to fairness on different GAM algorithms. Binder and Tutz [3] compared three different spline training algorithms, including backfitting, joint optimization, and boosting, finding that boosting performed particularly well in high-dimensional settings. Lou et al. [19] also found that boosted shallow bagged trees yielded higher accuracy than other GAM algorithms. However both papers focused on accuracy, not interpretability.

3 METHODS

In this section we describe the different GAM algorithms used in this paper. To make it easier for readers, we defer the description of the new metrics we define in this paper – feature sparsity and data fidelity – to just before their use in Sec. 5.

3.1 GAM Algorithms

Given an input $x \in \mathbb{R}^{N \times D}$, a label y , a link function g (e.g. in binary classification, g is logit), and shape functions f_j for each feature, a generalized additive models (GAM) can be written as:

$$g(y) = f_0 + \sum_{j=1}^D f_j(x_j). \quad (1)$$

GAMs are interpretable because the impact of each feature, f_j , on the prediction can be visualized as a graph (see Fig. 2 for an example), and humans can easily simulate how a GAM works by reading f_j s off different features from the graph and adding them together. We select the following six GAM algorithms to compare in this paper based on their popularity, state-of-the-art performance and availability of open source implementations.

Explainable Boosting Machine (EBM) A tree-based GAM designed for intelligibility and high accuracy [4, 19, 23] where shape functions f_j are gradient-boosted ensembles of bagged trees. Each tree operates on a single variable, preventing interactions effects from being learned. Trees are grown by repeatedly cycling through features, which forces the model to sequentially consider each feature as an explanation of the current residual rather than greedily selecting the best feature. This deliberate construction makes this model have less *feature sparsity*. For comparison, we create a sparse version of EBM similar to regular gradient boosted trees, “EBM-BF” (EBM-BestFirst), that greedily grows the next tree on the best, most informative feature to reduce error as much as possible at each step. Like most gradient boosted trees, EBM-BF is likely to put most weight on a few very important features, modest weight on a larger number of moderately useful features, and little

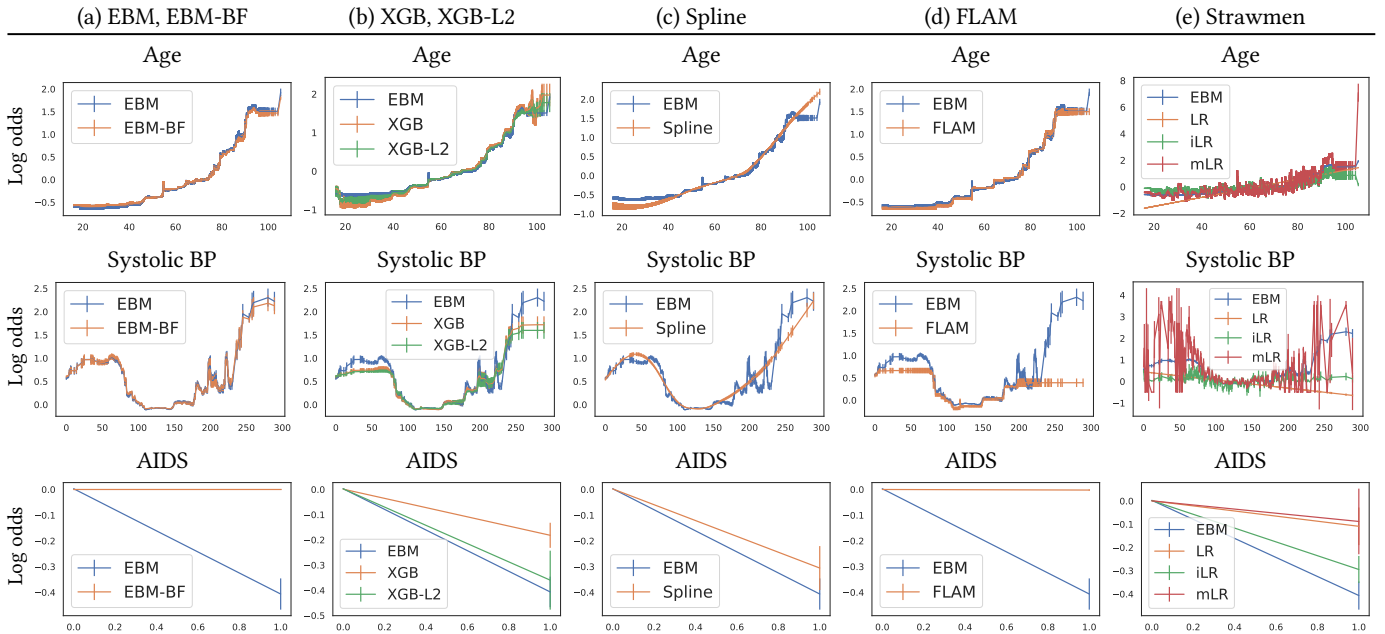


Figure 2: Shape plots from nine GAM algorithms trained on the MIMIC-II dataset (three of seventeen features shown). To make comparisons easier EBM (blue line) is repeated in each plot.

or no weight on features whose signal could be learned by other stronger, correlated features.

XGBoost (XGB) We introduce a new tree-based GAM based on the popular boosting package XGBoost [5]. To convert XGB to a GAM, we limit tree depth to 1 (stumps) so that the trees are not able to learn feature interactions, and we bag XGB to improve accuracy (similar to EBM). We also create a new version of XGB, "XGB-L2", similar to EBMs, that picks features sequentially when growing trees instead of greedily choosing the best feature. To achieve this, we set the XGB random subsampling of features parameter to a small ratio such that each tree is given just 1 feature. This deliberate modification makes this model have less feature sparsity. This modification makes XGB more of a "dense" model similar to ℓ_2 regularization that often uses all features. Fig. 2(b) shows these 2 methods. To our surprise, although XGB and EBM are both boosted trees, their shape plots can be quite different (Fig. 2(b)).

Spline A classic way to train GAMs is with spline basis functions [10]. We tried a variety of spline methods in 2 popular packages, the Python pygam [32] and R mgcv package [42], and chose cubic splines in pygam because it has a good combination of accuracy, robustness and speed (Fig. 2(c)).

Fused LASSO Additive Models (FLAM) For each unique value of feature x_j , Fused LASSO Additive Model (FLAM) [25] learns a weight on each value, and adds an ℓ_1 penalty to the differences between adjacent weights. This ℓ_1 penalty causes FLAM to produce relatively flat graphs and penalize unnecessary jumps. We use the R package FLAM [25] in our experiments (Fig. 2(d)).

Logistic regression (LR) and other strawmen approaches We compare these other approaches to Logistic Regression (LR), a widely used linear model that cannot learn non-linear shape plots.

We also compare to two other strawmen approaches: marginalized LR (**mLR**) and indicator LR (**iLR**). We first bin each feature x_j into at most 255 bins. In contrast to LR that assumes $f_j(x_j) = w_j x_j$, **mLR** sets $f_j(x_j) = w_j g(x_j)$ where $g(x_j)$ is the average (marginalized) value of target y within the same bin as x_j in the dataset. This is a GAM model built by applying logistic regression on top of marginalization, thus preventing shape plots from being learned in concert with each other. **iLR** treats each bin as a new feature (similar to one-hot encoding) and learns an LR on the transformed features. It thus ignores proximity relationships between different feature values (Fig. 2(e)).

3.2 Training and Hyperparameters

To fairly compare different GAM algorithms, we choose hyperparameters that perform best for each algorithm individually. Below, we briefly mention how we tune each GAM algorithm, and point the reader to the Appendix for more details.

For tree-based methods **EBM** and **XGB**, we perform early stopping to determine the optimal number of trees, stopping when the validation set performance stops improving for more than 50 trees. For **Spline**, we choose a maximum of 50 knots and use the gcv criterion [40] to select the smoothness penalty. We found that using more than 50 knots is intractable for larger datasets and does not improve performance in smaller datasets. For **FLAM**, we cross-validate the λ parameter and then refit the model on the entire training set using the optimal parameter. For **LR**, we cross-validate the ℓ_2 penalization parameter.

We split each dataset into 70-15-15% train-val-test splits and repeat our training procedure run 5 times. This allowed us to derive

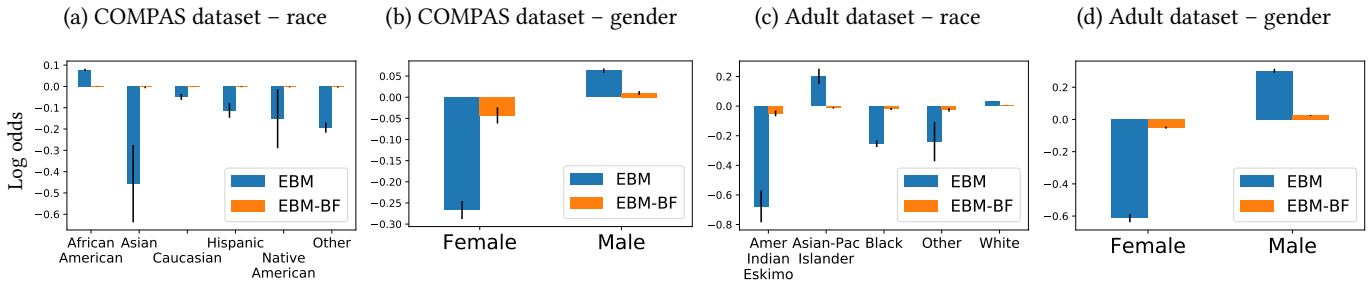


Figure 3: Shape plots for COMPAS and Adult datasets for two sensitive attributes: race and gender. We compare two extreme GAMs: a dense-feature GAM (EBM) and a sparse-feature GAM (EBM-BF). Sparse EBM-BF learns little effect on these features.

uncertainty estimates in the form of standard deviation across multiple runs.

4 CASE STUDIES: COMPAS, ADULT, MIMIC-II

We start with some case studies to highlight the implications of different GAM algorithms on common interpretability tasks such as surfacing unfairness or discovering anomalies in data. In this section, we highlight our key findings with plots specifically picked to be representative of our main results. A complete set of plots can be found in Appendix B.

4.1 How feature sparsity affects fairness?

One key property we study in this paper is which GAM algorithm uses fewer features to make predictions i.e. *feature sparsity*. Although sparsity is sometimes preferred because it appears to generate simpler explanations, it can hide data bias and discriminate against minority groups. Here we examine the sparsity properties of different GAM algorithms on two datasets that have been studied in the fairness community for racial and gender bias [6, 21, 44]. The COMPAS dataset contains demographic, crime, and recidivism information for defendants in Broward County, Florida, in 2013 and 2014. Research has suggested that the COMPAS recidivism risk score may be racially biased [1]. The Adult dataset extracted demographic information, including age, race, occupation, sex, etc. from the 1994 census data to predict if an individual’s income exceeds 50k/yr. In the dataset, males have on average higher annual incomes than females [21].

To motivate our analysis, we compare two GAM algorithms that are very different from each other in terms of feature sparsity: sparse EBM-BF and regular, “dense” EBM, yet achieve similar accuracy (see Table 4). Figure 3 displays the shape plots on two sensitive attributes, race and gender, on the COMPAS dataset. Since these features have modest influence compared to other features, the sparse-feature EBM-BF shows no or only a tiny effect on these sensitive attributes, while EBM shows much larger effects. Although there is no easy way to judge which GAM is more “causally” correct, the sparse EBM-BF makes users unaware of bias that may exist in the data and has been learned by other stronger, correlated features. In contrast, the dense EBM shows effects on all features. Because of this, we suggest that the dense model is better suited for surfacing potential bias in data than can then be investigated further by humans.

Next, we investigate how feature sparsity affects minority groups. Table 1 presents the predictive performance (cross entropy loss) of EBM and EBM-BF on each minority group. Although EBM and EBM-BF have negligible difference (less than 0.5%) in terms of overall loss, compared to EBM, EBM-BF exhibits greater loss on minority groups Other (1.45%) and Asian (6%) compared to majority group White (−0.02%); EBM exhibits lower loss on the Native American group (−2.26%). To further investigate this phenomenon, we perform an ablation study by removing the race feature from EBM thus forcing EBM to be more sparse. While this increased overall loss by 0.1% compared to EBM with the race feature, the loss for minority groups was again substantially increased, with the loss increasing by 6% for Asian and 1% for Native American. Similarly, when we remove the sex feature from EBM, the loss for the minority group Female increased by 0.99%, almost four times larger than the overall loss increase (0.23%). Unexpectedly, removing the sex feature improves the loss for minority group Native Americans (−5.32%); this is a possible explanation for why the loss for Native Americans is smaller for EBM-BF than EBM, as EBM-BF placed little importance on sex.

We repeat the same analysis on the Adult dataset. Table 2 presents the loss of EBM and EBM-BF on each minority group in the Adult dataset. Compared to EBM, EBM-BF exhibit greater loss on minority groups Indian (7.13%) and Other (19.05%), much more than the overall loss (5.27%) or loss on majority group White (5.17%). We also find that removing race from the EBM model increased the loss more for minority groups Indian (5.61%) and Other (1.04%), and removing sex from the EBM model increases the loss for Female (5.78%) much more than for Male (1.54%).

Implications GAM algorithms with a tendency to use fewer features to make predictions (e.g. EBM-BF) showed only small effects on sensitive attributes and exhibited greater prediction loss on minority groups causing unfairness, compared to GAM algorithms that tend to use more features to make predictions (e.g. EBM).

4.2 Data Anomaly Discovery

Another key property we study in this paper is which GAM algorithm is better able to capture anomalies in data. To illustrate, we train different GAM algorithms on a medical dataset: ICU mortality prediction dataset MIMIC-II [17]. On this dataset, XGB and EBM have similar shape plots thus we only present the EBM plots here for simplicity.

Table 1: Cross entropy loss of GAMs on different subpopulations in the COMPAS dataset. n is the number of samples in the subpopulation. The percentage shown is relative to the performance of EBM. Columns are sorted by descending n.

	All (n=6172)	Black (n=3175)	White (n=2103)	Other (n=343)	Asian (n=31)	Native American (n=11)	Male (n=4997)	Female (n=1175)
EBM	0.586	0.591	0.590	0.542	0.470	0.571	0.591	0.564
EBM-BF	0.589 (0.49%)	0.595 (0.72%)	0.590 (-0.02%)	0.550 (1.45%)	0.500 (6.48%)	0.558 (-2.26%)	0.594 (0.41%)	0.569 (0.82%)
EBM without race	0.587 (0.10%)	0.591 (0.06%)	0.590 (-0.01%)	0.544 (0.31%)	0.498 (6.08%)	0.579 (1.39%)	0.593 (0.18%)	0.563 (-0.30%)
EBM without sex	0.588 (0.23%)	0.594 (0.57%)	0.588 (-0.23%)	0.547 (0.95%)	0.464 (-1.14%)	0.540 (-5.32%)	0.592 (0.06%)	0.570 (0.99%)

Table 2: Cross entropy loss of GAMs on different subpopulations in the Adult dataset.

	All (n=32561)	White (n=27816)	Black (n=3124)	Asian/Pac (n=1039)	Indian/Eskimo (n=311)	Other (n=271)	Male (n=21790)	Female (n=10771)
EBM	0.265	0.277	0.163	0.309	0.204	0.137	0.321	0.152
EBM-BF	0.279 (5.27%)	0.291 (5.17%)	0.171 (5.37%)	0.326 (5.61%)	0.219 (7.13%)	0.164 (19.05%)	0.336 (4.78%)	0.163 (7.39%)
EBM without race	0.265 (0.15%)	0.277 (0.02%)	0.164 (0.98%)	0.311 (0.73%)	0.216 (5.61%)	0.139 (1.04%)	0.321 (0.12%)	0.152 (0.26%)
EBM without sex	0.271 (2.34%)	0.283 (2.27%)	0.170 (4.77%)	0.313 (1.29%)	0.201 (-1.62%)	0.138 (0.47%)	0.326 (1.54%)	0.161 (5.78%)

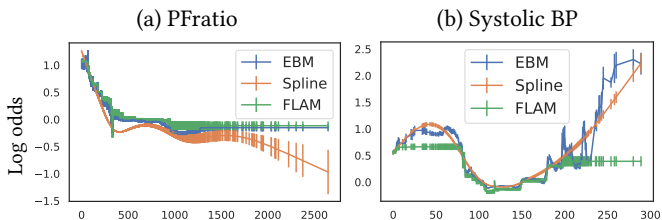


Figure 4: Two data anomalies in MIMIC-II that can be detected by tree-based GAMs (EBM): (a) PFratio missing values imputed using population mean 332; (b) Systolic BP with likely human intervention artifacts at 175, 200 and 225.

Fig. 4(a) displays one feature, PFratio (a measure of how well patients convert oxygen in air to oxygen in blood), for the three most accurate GAM algorithms on this dataset: EBM, Spline and FLAM. Interestingly, both EBM and FLAM capture a sharp drop in mortality risk at PFratio=332. It turns out that PFratio is usually not measured for healthier patients, and the missing values for these patients have been imputed by its population mean 332 (a common preprocessing fix for missing data), thus giving a group of low-risk patients the mean value of this feature. However, Spline is unable to represent the sharp drop, becoming distorted in the region 300-600, thereby underestimating the risk for patients in this region.

Fig. 4(b) for Systolic Blood Pressure (BP) shows another data anomaly that is only captured by tree-based GAM algorithms EBM

and XGB. EBM captures three jumps, exhibiting dips in risk predictions near 175, 200 and 225. These are likely to be human intervention artifacts, since 175, 200, and 225 are treatment thresholds used by physicians. As a patient’s Systolic BP increases the mortality risk naturally increases, but when they reach the next treatment threshold, risk actually drops because most patients just above the threshold are receiving more aggressive care that is effective at reducing their risk. Both Spline and FLAM are too smooth or flat and fail to capture these anomalies.

Implications Localized data anomalies such as mean imputation and human intervention artifacts (e.g. medical treatment thresholds), often require models to learn quick, non-linear changes in risk. Tree-based methods (e.g. EBM and XGB) can detect these much better compared to GAM algorithms that are too smooth or sparse (e.g. Spline and FLAM).

5 QUANTITATIVE ANALYSIS OF GAMs

In the previous section, we saw examples of how different GAM algorithms revealed different insights. In this section, we study the performance differences between GAM algorithms quantitatively. We first benchmark the test accuracy of different GAMs on ten different datasets (Sec. 5.1). Then we measure feature sparsity of different GAM algorithms (Sec. 5.2). Finally, we measure data fidelity using both real (Sec. 5.3) and simulated data (Sec. 5.4, 5.5).

Table 3: Description of ten classification datasets used.

	Domain	N	P	Positive Rate	Description
Adult	Finance	32,561	14	24.08%	Income prediction
Breast	Healthcare	569	30	62.74%	Cancer classification
Churn	Retail	7,043	19	26.54%	Subscription churning
Credit	Retail	284,807	30	0.17%	Fraud detection
COMPAS	Law	6,172	6	45.51%	Reoffense risk scores
Heart	Healthcare	457	11	45.95%	Heart Disease
MIMIC-II	Healthcare	24,508	17	12.25%	ICU mortality
MIMIC-III	Healthcare	27,348	57	9.84%	ICU mortality
Pneumonia	Healthcare	14,199	46	10.86%	Mortality
Support2	Healthcare	9,105	29	25.92%	Hospital mortality

5.1 GAM accuracy

How do we choose which GAM to use? Accuracy is perhaps the first obvious consideration. Table 4 provides test set AUC of different GAM algorithms on ten datasets. These datasets of varying size (500 - 250k samples) and number of features (6 - 57 features) span different domains such as healthcare, criminal justice, finance, and retail (see Table 3). In addition to the nine GAM algorithms described in Sec. 3, we also include two full-complexity methods: Random Forest (RF) and XGB with depth 3 (XGB-d3). For each method, we compute three metrics, each of which is averaged over ten datasets: (1) Test AUC; (2) Rank of test AUC compared to other methods (lower rank is better); (3) Test AUC normalized compared to other methods (lowest test AUC for a dataset has value 0, highest test AUC for a dataset has value 1, with all other test AUCs scaled linearly between them). On average across ten datasets, EBM, EBM-BF, and XGB-d3 performed the best. In general, GAMs perform better than or comparably to full complexity models. Four of the GAMs (EBM, XGB, Spline and FLAM) achieve similar top performance with average AUC differences less than 0.2%.

Implications There exist GAM algorithms that perform comparably to full complexity models. Several GAM algorithms are similarly accurate, hence accuracy should not be the sole consideration when selecting between different GAM algorithms.

5.2 GAM feature sparsity

In this section we propose a new metric to quantify feature sparsity, the notion that some GAM algorithms use fewer features than others to make predictions, which we have seen in Sec. 4.1 to impact bias discovery.

Feature density metric The idea is to quantify how fast the test *error* of a trained model decays (i.e., how fast the model becomes more accurate) as we allow the model to have access to more features; a sparse model only requires a few important features to quickly reduce its test error, while a dense model needs more features to recover because it will have spread learned effects across more of the features. Using the GAM formulation as in Equation 1, we proceed as follows to compute this metric: first we keep only f_0 and measure the GAM’s test set error as the initial error. Then for each step out of D steps, we greedily search over which feature $f_j(x_j)$, when added back to the model, reduces its validation error the most. We add that feature back and measure how the model’s

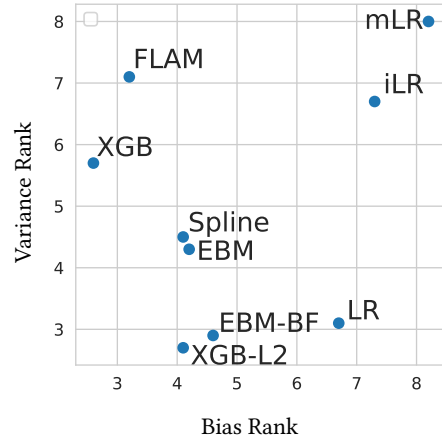


Figure 5: Bias rank (x-axis) vs. variance rank (y-axis) for each GAM across multiple datasets. Lower rank is better.

test error decreases. We save the test error as each subsequent feature is added, until D features are added after D steps, and plot test error against features. Finally we compute the feature density metric as the normalized area under this curve, treating the initial test error as 100 and final error (with D features) as 0. We expect an extremely sparse model to have value close to 0, and a dense model to have value close to 50.

Table 5 presents normalized feature density for different GAM algorithms on ten datasets. As expected, EBM consistently has higher density (less sparsity) than EBM-BF across datasets, as it uses more features by design. Similarly, XGB-L2 also has higher density than XGB, and LR is higher than LASSO. This confirms that the feature density metric reflects what we want. FLAM has low feature density, which is unsurprising due to the ℓ_1 penalty present in the method. Spline does not exhibit a clear pattern of feature density. For example, Spline has the smallest density on the Adult dataset but the largest density on the Breast dataset.

Implications The proposed feature density metric captures expected behavior. We see lower feature density for methods that greedily select the next best feature (e.g. EBM-BF) or have penalties that regularize for sparsity (e.g. FLAM). Methods that repeatedly cycle over all features (e.g. EBM) have higher feature density.

5.3 GAM data fidelity

In this section we propose a new metric to quantify how well a GAM is able to capture underlying data patterns, which we have seen in Sec. 4.2 to impact data anomaly discovery.

At first glance, one may think that test accuracy is a suitable metric for this purpose, since it captures how well a model generalizes to unseen data. However, we saw in Sec. 4 when comparing GAM algorithms of similar test accuracy how some were less able to represent certain data patterns. For example, smooth basis functions in Spline, while reducing variance and hopefully improving test set generalization, limited the model’s ability to capture sharp jumps in the data. As noted by Shmueli [33], some highly accurate predictive models may actually be “wrong” in terms of capturing underlying

Table 4: Test set AUCs (%) across ten datasets average over five runs. Best number in each row is in bold.

	GAM									Full Complexity	
	EBM	EBM-BF	XGB	XGB-L2	FLAM	Spline	iLR	LR	mLR	RF	XGB-d3
Adult	93.0 ± 0.5	92.8 ± 0.5	92.8 ± 0.6	91.7 ± 0.6	92.5 ± 0.6	92.0 ± 0.6	92.7 ± 0.5	90.9 ± 0.6	92.5 ± 0.4	91.2 ± 0.5	93.0 ± 0.6
Breast	99.7 ± 0.5	99.5 ± 0.5	99.7 ± 0.5	99.7 ± 0.5	99.8 ± 0.3	98.9 ± 0.8	98.1 ± 0.5	99.7 ± 0.4	98.5 ± 0.5	99.3 ± 1.1	99.3 ± 1.1
Churn	84.4 ± 0.7	84.0 ± 0.9	84.3 ± 0.7	84.3 ± 0.7	84.2 ± 0.7	84.4 ± 0.8	83.4 ± 1.0	84.3 ± 0.7	82.7 ± 1.0	82.1 ± 0.6	84.3 ± 0.7
COMPAS	74.3 ± 1.4	74.5 ± 1.7	74.5 ± 1.5	74.3 ± 1.5	74.2 ± 1.7	74.3 ± 1.5	73.5 ± 1.3	72.7 ± 1.0	72.2 ± 1.3	67.4 ± 1.2	74.5 ± 1.5
Credit	98.0 ± 0.5	97.3 ± 1.3	98.0 ± 0.6	98.1 ± 0.6	96.9 ± 0.4	98.2 ± 0.7	95.6 ± 0.6	96.4 ± 1.1	94.0 ± 1.4	96.2 ± 1.5	97.3 ± 0.7
Heart	85.5 ± 6.9	83.8 ± 6.0	85.3 ± 6.3	85.8 ± 7.0	85.6 ± 6.7	86.7 ± 6.3	85.9 ± 6.3	86.9 ± 5.8	74.4 ± 5.3	85.4 ± 6.5	84.3 ± 4.6
MIMIC-II	83.4 ± 0.9	83.3 ± 0.8	83.5 ± 1.0	83.4 ± 0.9	83.4 ± 1.0	82.8 ± 0.8	81.1 ± 1.0	79.3 ± 0.8	81.6 ± 0.7	86.0 ± 0.6	84.7 ± 0.7
MIMIC-III	81.2 ± 0.4	80.7 ± 0.7	81.5 ± 0.5	81.5 ± 0.5	81.2 ± 0.4	81.4 ± 0.4	77.4 ± 1.0	78.5 ± 0.5	77.6 ± 0.3	80.7 ± 0.8	82.0 ± 0.7
Pneumonia	85.3 ± 0.6	84.7 ± 0.7	85.0 ± 0.8	85.0 ± 0.6	85.3 ± 0.9	85.2 ± 0.6	84.3 ± 1.0	83.7 ± 0.6	84.5 ± 0.7	84.5 ± 0.5	84.8 ± 0.8
Support2	81.3 ± 1.0	81.2 ± 1.0	81.4 ± 1.1	81.2 ± 1.0	81.2 ± 1.1	81.2 ± 1.1	80.0 ± 1.2	80.3 ± 0.7	77.2 ± 0.9	82.4 ± 1.0	82.0 ± 1.4
Average AUC	86.6	86.2	86.6	86.5	86.4	86.5	85.2	85.3	83.5	85.5	86.6
Average Rank	3.70	6.70	3.40	4.90	5.05	4.60	8.70	7.75	9.70	7.40	4.10
Normalized AUC	89.3	78.1	87.3	81.8	83.6	81.0	47.4	50.7	28.5	54.3	86.5

Table 5: Normalized feature density (%). Higher numbers mean the model uses more features. Highest number in each row is in red; lowest number in each row is in blue. Columns are sorted by descending density.

	XGB-L2	EBM	LR	Spline	XGB	LASSO	FLAM	EBM-BF
Adult	33.9	27.1	22.0	20.5	29.0	21.3	21.1	22.6
Breast	11.2	08.6	13.0	23.4	07.0	06.6	07.7	05.9
Churn	15.0	15.7	19.9	22.7	12.9	16.2	13.1	13.0
COMPAS	18.3	18.3	17.7	17.3	17.9	17.7	17.2	17.0
Credit	26.9	26.9	12.4	19.1	19.4	12.2	17.0	15.8
Heart	28.7	24.0	32.6	15.4	25.0	30.8	21.5	21.8
MIMIC-II	20.5	20.4	19.4	21.0	19.6	19.4	18.8	18.6
MIMIC-III	21.2	20.7	19.0	21.6	18.7	18.7	18.6	14.8
Pneumonia	29.9	29.7	27.2	19.5	25.3	25.8	25.8	20.6
Support2	11.4	12.4	10.3	12.6	13.0	10.2	11.4	11.7
Average	21.7	20.4	19.4	19.3	18.8	17.9	17.2	16.2

data patterns. This notion is exactly statistical bias, which arises from model misspecification of the underlying data patterns [12].

Data fidelity metric We use an approximation to the bias term in a bias-variance analysis to measure data fidelity. In bias-variance analysis [2], the loss of model is composed of noise $N(x)$, bias $B(x)$ and variance $V(x)$ terms:

$$E_{D,t}[L(t, y)] = N(x) + B(x) + V(x) \quad \text{where}$$

$$N(x) = E_t[L(t, y_*)], B(x) = L(y_*, y_m), V(x) = E_D[L(y_m, y)]$$

where D is the training distribution, t is the true label, y_* is the optimal predictions, y_m is the mean prediction of models across possible training datasets, and y is the model. Since we do not know the y_* , we instead measure the *empirical bias* combining both noise and bias $N(x) + B(x) = E_t[L(t, y_m)]$ following Munson and Caruana [22]. We use the following sampling procedure: in each round, we split our dataset into 85-15% train-test splits. We then randomly subsample the training data to 50% and train models 5 times, and we set the average of 5 models as y_m to calculate

empirical bias and variance once. Finally, the bias and variance estimates are averaged over eight rounds, and ranked compared to other GAM algorithms on each dataset. We take the average ranks across the ten datasets (lower rank is better).

Fig. 5 plots average variance rank vs. average bias rank for different GAM algorithms. Considering GAM algorithms closest to the bottom left corner (i.e. (0, 0) point), which are also the most accurate GAMs (see Table 4), XGB has the highest data fidelity (lowest bias rank) but has rather high variance. FLAM has the next highest data fidelity, but has even higher variance, hence it is dominated by XGB that has both higher data fidelity and lower variance. After FLAM, XGB-L2, Spline, and EBM have the next highest data fidelity, and promisingly, with significantly lower variance than FLAM or XGB.

Implications We use statistical bias as a proxy to measure data fidelity with real data. By decomposing error into bias and variance components, we see that equally accurate GAM algorithms achieve the same accuracy in different ways. Certain GAM algorithms (e.g. XGB) have lower bias which indicates better fidelity, while other GAMs (e.g. XGB-L2) have lower variance at the expense of higher bias.

5.4 GAM data fidelity and generator bias

We have thus far studied the data fidelity properties of different GAM algorithms on several real datasets. However, it may be that the inductive bias of a certain GAM algorithm happened to agree with the (unknown) data pattern in a particular real dataset. In this section, we experiment with semi-synthetic datasets created using known data generators. To preserve the character of real datasets as much as possible, we keep the features X but change the label y by training multiple ground truth GAM models (EBM, XGB, Spline, FLAM and LR) on features X and then *re-generating* the label y as each model’s predictions. Since these GAM models (except LR) are among the most accurate models on most datasets (Table 4), the generated labels capture the real-world distribution as close as possible. As these GAM algorithms are very different from each other, this should provide a diversity of ground truth data patterns.

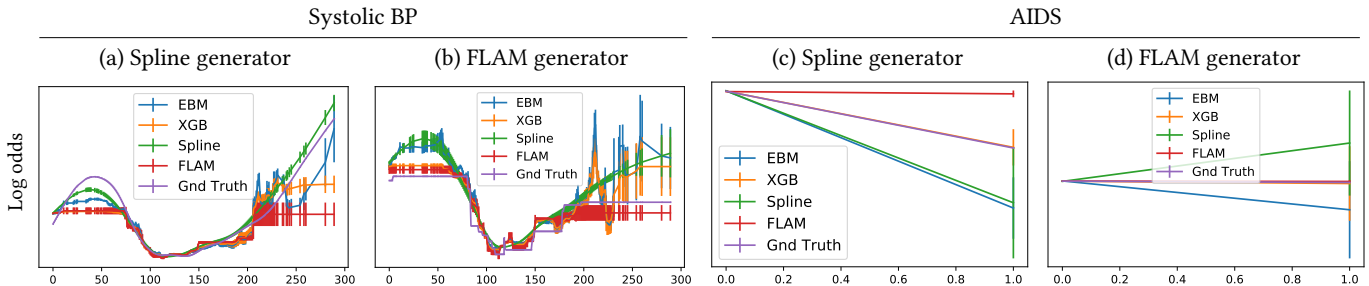


Figure 6: Shape plots for Systolic BP and AIDS features in semi-synthetic MIMIC-II, generated using different generators (Spline and FLAM).

Fig. 6(a)-(d) shows different GAMs alongside ground truth patterns from two very different generators, Spline and FLAM, on MIMIC-II for one continuous feature (Systolic BP) and one boolean feature (AIDS). Purple represents ground truth, i.e. Spline generator for Fig. 6(a) and (c), and FLAM generator for 6(b) and (d). We see an obvious *generator bias*: a GAM algorithm fits the ground truth better when ground truth is generated using the same algorithm. For example, on Systolic BP, the Spline GAM fits well the data generated by its own generator (Fig. 6(a)), while doing poorly for data generated by the FLAM generator (Fig. 6(b)), and vice versa for FLAM. However, tree-based methods (EBM, XGB) on Systolic BP with the Spline generator (Fig. 6a) still learn abrupt jumps at 225 even when the underlying ground truth is smooth; similarly there is also a drop at 175. This illustrates that it is possible for model inductive bias to dominate irrespective of the true data generator.

To mitigate the aforementioned generator bias, we perform a worst-case analysis: what is the worst performance each GAM algorithm would get across all of the different data generators? Since we do not know the underlying generators on real datasets – they could be jumpy, smooth, or even linear – this analysis is more realistic and robust to all these cases.

Worst-case data fidelity metric To measure how well a GAM can recover the ground truth generators, we calculate the mean absolute difference of each shape plot between the ground truth GAM and the GAM model. Specifically, using the GAM formulation as in Equation 1 where f_j is the shape function for feature j , and taking g_j to be the shape function for the ground truth GAM, we calculate the absolute difference $\sum_{j=1}^D |f_j(x_j) - g_j(x_j)|$ across the whole dataset. To compare between datasets, we linearly scale the absolute difference between 0 and 100 for a particular semi-synthetic dataset, with the worst GAM algorithm having value 0 and best GAM algorithm having value 100. We then take the worst score over the five different data generators that yielded five semi-synthetic datasets from each real dataset.

Table 6 provides the worst-case data fidelity for eight GAM algorithms on six real datasets, where each dataset (row) encapsulates five semi-synthetic datasets from different data generators. FLAM and XGB performed the best, then EBM and Spline.

Table 6: Worst-case data fidelity (%) taking into account different data generators. Each row aggregates the results over five different generators (EBM, XGB, FLAM, Spline and LR). Higher numbers are better. Best number in each row is in bold.

	FLAM	XGB	EBM	Spline	EBM-BF	LR	iLR	mLR
Breast	22.9	30.3	0	13.3	21.2	42.5	0	0
Churn	20.3	10.5	13.5	0	1.3	0	16.0	0
Heart	86.9	68.2	68.7	69.7	24.8	52.4	64.6	0
MIMIC-II	62.9	73.9	61.2	72.7	52.6	0	6.6	0
MIMIC-III	65.2	70.1	45.3	51.2	37.0	27.0	0	0
Pneumonia	64.4	60.2	40.0	3.6	0	0	26.8	6.4
Average	53.8	52.2	38.1	35.1	22.8	20.3	19.0	1.1

Implications FLAM and XGB exhibit the best worst-case data fidelity. Spline and EBM are similar, and EBM-BF is the worst. Taking into account different data generators, our results are not substantively different from the results derived from the bias-variance analysis on real data in Sec. 5.3.

5.5 GAM accuracy vs. data fidelity

A GAM model that has high accuracy but low data fidelity may mislead users who tend to judge models solely based on accuracy. We quantify which GAM algorithm is more likely to mislead users this way, by comparing the difference between test AUC rank and data fidelity rank. For each dataset, we compute these two ranks as in Sec. 5.1 and Sec. 5.3, with lower rank being better. Then we take the rank of fidelity minus the rank of test AUC. If the result is negative, we clip it at 0. We call this the “positive difference” between the two ranks. Finally, we average this over all thirty semi-synthetic datasets. We expect a misleading model to have a lower test AUC rank and higher data fidelity rank.

From Table 7, Spline has the largest difference in rank over multiple datasets with different data generators. This rank difference is largest when the data generators are jumpy, which creates challenges for Spline which uses smooth basis functions.

Implications For Spline, using high test accuracy alone to select a model may be misleading, especially when the underlying data pattern may be jumpy. Other methods are more stable.

Table 7: Difference between test AUC rank and data fidelity rank on thirty semi-synthetic datasets. The larger this difference, the less reliable it is to use high accuracy to infer good data fidelity. Best number is bold.

	EBM	FLAM	XGB	LR	EBM-BF	Spline
Avg Pos Diff in Rank	0.47	0.50	0.62	0.63	0.87	1.22

Table 8: Summary of key findings, ranking the different GAM algorithms across six properties studied in this paper. Best number in each row is in bold or red.

	EBM	XGB	FLAM	Spline	LR
Test-set accuracy	1.5	1.5	4	3	5
Feature density	1	4	5	2.5	2.5
Low bias (bias/variance)	3.5	1	2	3.5	5
Worst-case fidelity	3.5	1.5	1.5	3.5	5
High accuracy implies good data fidelity	1.5	3.5	1.5	5	3.5
Anomaly detection	1.5	1.5	3	4.5	4.5
Sum of Ranks	12.5	13	17	22	25.5

6 DISCUSSION

GAMs are widely used to discover patterns in data in a variety of fields including business [30], healthcare [15], ecology [24], horticulture [31], air pollution [26], nutrition [28] and COVID-19 [16]. But most of these research only experimented with a specific GAM algorithm (typically Spline) without any comparison to other GAM algorithms. In this work, we have shown that the patterns learned by GAMs are highly impacted by their own inductive biases. If the papers that used GAMs to discover patterns had used different GAM algorithms, would they have drawn different conclusions? How many of the findings are due to true patterns in the data and not due to the inductive bias of the particular GAM algorithm chosen?

While we aimed to provide a useful and fair experimental study, there are limitations to the conclusions that can be drawn from our work due to design choices we made. In terms of data sets, we considered common Kaggle datasets across several domains that are relatively large but still have a manageable amount of features. We do not explore small datasets used in the Spline literature, where a smoothing prior might help compensate for a lack of sample size. In terms of models, we only focused on a few of the most representative GAM algorithms and make additional modifications to these methods to study different characteristics of GAMs (e.g. feature sparsity and data fidelity). We leave more theoretical comparisons to future work.

7 CONCLUSION

The key findings are summarized in Table 8, where we have synthesized our findings across six different properties studied in this paper and ranked each GAM algorithm for each property (ties count for half a rank). Although a number of GAM algorithms yield similar accuracy, tree-based methods like EBM and XGB are superior when considering issues such as bias and data anomaly discovery, sparsity, fidelity, and accuracy. Tree-based methods such as XGB and EBM have higher feature density than FLAM or Spline. They also have less bias on real data, and recover data patterns with better fidelity on semi-synthetic data. We also find Spline could have high accuracy yet at the same time low data fidelity, which might mislead users who perform model selection based on test accuracy alone. Qualitatively, Spline and FLAM are not good at detecting local anomalies such as mean imputation or treatment effects, both of which are easily detected by the tree-based methods. Spline also extrapolates over-confidently in low-sample regions (Fig. 1(b), and see other examples in Appendix B).

Future development of better GAM algorithms should focus on the following: (1) GAMs that can better capture rapid non-linear change, (2) GAMs with high feature density to improve fairness and prevent bias masking, (3) GAMs having higher data fidelity on both real and simulated data. We believe our work is an important step towards making GAMs more trustworthy, and our evaluation framework will promote the development of better GAMs in the future.

ACKNOWLEDGMENTS

This work was created during an internship at Microsoft Research. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute (www.vectorinstitute.ai/#partners).

REFERENCES

- [1] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2019. Machine Bias: There’s software used across the country to predict future criminals. And it’s biased against blacks. (2019). <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- [2] Eric Bauer and Ron Kohavi. 1999. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine learning* 36, 1 (1999).
- [3] Harald Binder and Gerhard Tutz. 2008. A comparison of methods for the fitting of generalized additive models. *Statistics and Computing* 18, 1 (2008).
- [4] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *KDD*.
- [5] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *KDD*.
- [6] Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data* 5, 2 (2017).
- [7] Finale Doshi-Velez and Been Kim. 2017. Towards A Rigorous Science of Interpretable Machine Learning. *Springer Series on Challenges in Machine Learning: "Explainable and Interpretable Models in Computer Vision and Machine Learning"* (2017).
- [8] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [9] Kenneth Dwyer and Robert Holte. 2007. Decision Tree Instability and Active Learning. In *ECML*.
- [10] Trevor Hastie and Rob Tibshirani. 1990. *Generalized Additive Models*. Chapman and Hall/CRC.
- [11] Trevor Hastie and Robert Tibshirani. 1995. Generalized additive models for medical research. *Statistical Methods in Medical Research* 4, 3 (1995).
- [12] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science &

- Business Media.
- [13] Stefan Hegselmann, Thomas Volkert, Hendrik Ohlenburg, Antje Gottschalk, Martin Dugas, and Christian Ertemer. 2020. An Evaluation of the Doctor-Interpretability of Generalized Additive Models with Interactions. In *Machine Learning for Healthcare Conference*.
- [14] Giles Hooker and Lucas Mentch. 2019. Please Stop Permuting Features: An Explanation and Alternatives. *arXiv preprint arXiv:1905.03151* (2019).
- [15] Paul R Hunter and Annette Prüss-Ustün. 2016. Have we substantially underestimated the impact of improved sanitation coverage on child health? A generalized additive model panel analysis of global data on child mortality and malnutrition. *PLoS One* 11, 10 (2016).
- [16] Farzali Izadi. 2020. Generalized additive models to capture the death rates in Canada COVID-19. *arXiv preprint arXiv:1702.08608* (2020).
- [17] Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific Data* 3, 1 (2016).
- [18] Zachary C Lipton. 2018. The mythos of model interpretability. *Queue* 16, 3 (2018).
- [19] Yin Lou, Rich Caruana, and Johannes Gehrke. 2012. Intelligible models for classification and regression. In *KDD*.
- [20] Scott M Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *NeurIPS*.
- [21] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2019. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635* (2019).
- [22] M Arthur Munson and Rich Caruana. 2009. On feature selection, bias-variance, and bagging. In *ECML PKDD*.
- [23] Harsha Nori, Samuel Jenkins, Paul Koch, and Rich Caruana. 2019. InterpretML: A Unified Framework for Machine Learning Interpretability. *arXiv preprint arXiv:1909.09223* (2019).
- [24] Eric J Pedersen, David L Miller, Gavin L Simpson, and Noam Ross. 2019. Hierarchical generalized additive models in ecology: an introduction with mgcv. *PeerJ* 7 (2019).
- [25] Ashley Petersen, Daniela Witten, and Noah Simon. 2016. Fused lasso additive model. *Journal of Computational and Graphical Statistics* 25, 4 (2016).
- [26] Khaiwal Ravindra, Preeti Rattan, Suman Mor, and Ashutosh Nath Aggarwal. 2019. Generalized additive models: Building evidence of air pollution, climate change and human health. *Environment International* 132 (2019).
- [27] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. In *KDD*.
- [28] Maryam Rostami, Masoumeh Simbar, Mina Amiri, Razieh Bidhendi-Yarandi, Farhad Hosseinpanah, and Fahimeh Ramezani Tehrani. 2020. The optimal cut-off point of vitamin D for pregnancy outcomes using a generalized additive model. *Clinical Nutrition* (2020).
- [29] Veeranjaneyulu Sadhanala and Ryan J Tibshirani. 2019. Additive models with trend filtering. *The Annals of Statistics* (2019).
- [30] K Sapra. 2013. Generalized additive models in business and economics. *International Journal of Advanced Statistics and Probability* 1, 3 (2013).
- [31] Nay Min Min Thaw Saw, Claudio Moser, Stefan Martens, and Pietro Franceschi. 2017. Applying generalized additive models to unravel dynamic changes in anthocyanin biosynthesis in methyl jasmonate elicited grapevine (*Vitis vinifera* cv. Gamay) cell cultures. *Horticulture Research* 4, 1 (2017).
- [32] Daniel Servén and Charlie Brummitt. 2018. pyGAM: Generalized Additive Models in Python. <https://doi.org/10.5281/zenodo.1208723>
- [33] Galit Shmueli. 2010. To explain or to predict? *Statist. Sci.* 25, 3 (2010).
- [34] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. 2020. Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods. In *AIES*.
- [35] Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. 2007. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics* 8, 1 (2007).
- [36] Sarah Tan, Rich Caruana, Giles Hooker, Paul Koch, and Albert Gordo. 2018. Learning global additive explanations for neural nets using model distillation. *arXiv preprint arXiv:1801.08640* (2018).
- [37] Sarah Tan, Rich Caruana, Giles Hooker, and Yin Lou. 2018. Distill-and-compare: Auditing black-box models using transparent model distillation. In *AIES*.
- [38] Robert Tibshirani. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B* 58, 1 (1996).
- [39] Robert Tibshirani, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight. 2005. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B* 67, 1 (2005).
- [40] Grace Wahba. 1985. A comparison of GCV and GML for choosing the smoothing parameter in the generalized spline smoothing problem. *The Annals of Statistics* (1985).
- [41] Grace Wahba. 1990. *Spline models for observational data*. SIAM.
- [42] S. N. Wood. 2011. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society: Series B* 73, 1 (2011).
- [43] Marvin N Wright and Inke R König. 2019. Splitting on categorical predictors in random forests. *PeerJ* 7 (2019).
- [44] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *ICML*.
- [45] Zhengze Zhou and Giles Hooker. 2021. Unbiased Measurement of Feature Importance in Tree-Based Methods. *TKDD* (2021).

A REPRODUCIBILITY: TRAINING DETAILS, HYPERPARAMETERS, AND DATASETS

Code can be found at <https://github.com/zzzace2000/GAMs>.

A.1 Further training details and hyperparameters

In this section, we further describe training details and hyperparameters to supplement the discussion in Sec. 3.

- EBM, EBM-BF: we use the open-source package from <https://github.com/interpretml/interpret>. We set the parameters inner bagging as 100 and outer bagging as 100. We find that increasing the number of bags does not further improve performance. We use the default learning rate of 0.01, default early stopping patience set to 50, and the maximum 30000 episodes to make sure it converges.
- XGB, XGB-d3, XGB-L2: we use the open source package <https://xgboost.readthedocs.io/en/latest/index.html>. We also use the default learning rate with the same early stopping patience set as 50 and number of trees as maximum 30,000. We use bagging of 100 times and depth 1 for our XGB GAM model. For XGB-d3 (XGB with tree depth 3), we find that bagging of XGB-d3 hurts the performance a bit, and thus do not apply any bagging for XGB-d3. For XGB-L2, we set the parameter “colsample_bytree” as a small value 1e-5 to make sure each tree only sees one feature.
- FLAM: we use the package from R <https://cran.r-project.org/web/packages/flam/flam.pdf>. We use a 15% validation set to select the best λ penalty parameter in the fused LASSO, and then refit the whole data with the best penalty parameter. We set the parameter number of lambda as 100 and the minimum ratio as 1e-4 to increase the performance of the model.
- Spline: we use the pygam package [32]. We set the number of basis functions to be 50 and the maximum iteration as 500. We find increasing number of basis functions more than 50 would result in instability when fitting in large datasets.
- LR: we use scikit-learn’s LogisticRegressionCV with $C_s = 12$ (grid search for 12 different ℓ_2 penalty) and cross validation for 5 times to choose the best ℓ_2 , and re-fit on the whole data.
- iLR, mLR: we use the EBM package’s preprocessor to quantify bin the features into 255 bins. Then we use LR on top of it to train a linear model.

We also tried the following GAM algorithms but do not include them in the main results, for reasons detailed below:

- SKGBT: we try the gradient boosting tree in scikit-learn also with tree depth set as 1. The result is similar to EBM so we do not compare them in the paper.
- Cubic spline and plate spline in R mgcv package: to our surprise, mgcv is really unstable on two datasets, Breast Cancer and Churn. After some investigation, we found a possible reason to be that

mgcv does not handle numerical instability when the prediction is too close to 0 or 1.

A.2 Encoding categorical features

For datasets with categorical variables, the choice of encoding can affect both the shape plots and the accuracy. For gradient boosting trees, one may think that using label encoding (LE) is better than one-hot encoding, as one-hot encoding has been shown to have inferior performance in ensemble trees [43]. We investigate the effects of two types of encoding on EBM and XGB. In 6 of the datasets with categorical features, EBM with label encoding (LE) indeed shows superior performance to one-hot encoding. However, for XGB, one-hot encoding performs slightly better on average. Thus we use LE for EBM and one-hot encoding for XGB. For the rest of the methods, we use LE for mLR and one-hot encoding for FLAM, Spline, LR and iLR as these methods cannot handle inadequate numerical ordering.

A.3 Dataset sources

The datasets used in this paper can be found at:

- Adult: UCI [8]
- Breast cancer: UCI [8]
- Credit: <https://www.kaggle.com/mlg-ulb/creditcardfraud>
- Churn: <https://www.kaggle.com/blastchar/telco-customer-churn>
- COMPAS: <https://www.kaggle.com/danofer/compass>
- Heart disease: UCI [8]
- MIMIC-II and MIMIC-III dataset [17]
- Pneumonia: we thank the authors of Caruana et al. [4] for running our code on their dataset.
- Support2: <http://biostat.mc.vanderbilt.edu/DataSets>

B ADDITIONAL SHAPE PLOTS

The complete set of shape plots can be found at <https://drive.google.com/file/d/1PoMRgfuHYax6xuCVU0Dbut3yFJ2ohuLX/view?usp=sharing>.