
Learning Explainable Models Using Attribution Priors

Gabriel Erion*

Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98144
erion@cs.washington.edu

Joseph D. Janizek*

Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98144
jjanizek@cs.washington.edu

Pascal Sturmfels*

Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98144
pstorm@cs.washington.edu

Scott Lundberg

Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98144
slund1@cs.washington.edu

Su-In Lee

Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98144
suinlee@cs.washington.edu

Abstract

Two important topics in deep learning both involve incorporating humans into the modeling process: *Model priors* transfer information from humans to a model by constraining the model’s parameters; *Model attributions* transfer information from a model to humans by explaining the model’s behavior. We propose connecting these topics with *attribution priors*[†], which allow humans to use the common language of attributions to enforce prior expectations about a model’s behavior during training. We develop a differentiable axiomatic feature attribution method called *expected gradients* and show how to directly regularize these attributions during training. We demonstrate the broad applicability of attribution priors (Ω) by presenting three distinct examples that regularize models to behave more intuitively in three different domains: 1) on image data, Ω_{pixel} encourages models to have piecewise smooth attribution maps; 2) on gene expression data, Ω_{graph} encourages models to treat functionally related genes similarly; 3) on a health care dataset, Ω_{sparse} encourages models to rely on fewer features. In all three domains, attribution priors produce models with more intuitive behavior and better generalization performance by encoding constraints that would otherwise be very difficult to encode using standard model priors.

1 Introduction

Recent work on interpreting machine learning models has focused on *feature attribution methods*. Given an input feature, a model, and a prediction on a particular sample, such methods assign a

*Equal contribution (ordered alphabetically).

[†]<https://github.com/suinleelab/attributionpriors>

number to the input feature that represents how important the input feature was for making the prediction. Previous literature about such methods has focused on the axioms they should satisfy [18, 37, 35, 6], and how attribution methods can give us insight into model behavior [19, 20, 29, 41]. These methods can be an effective way of revealing problems in a model or a dataset. For example, a model may place too much importance on undesirable features, rely on many features when sparsity is desired, or be sensitive to high frequency noise (e.g. over pixels). In such cases, we often have a prior belief about how a model should treat input features, but for neural networks it can be difficult to mathematically encode this prior in terms of the original model parameters.

Ross et al. [26] introduce the idea of regularizing explanations in order to build models that both perform well and agree with domain knowledge. Given a binary variable indicating whether each feature should or should not be important for predicting on each sample in the dataset, their method penalizes input gradients of unimportant features. However, two drawbacks limit the method’s applicability to real-world problems. First, using gradients to determine feature importance suffers from problems like saturation and thresholding [33]. Second, it is often difficult to specify which features should be important in a binary manner.

More recent work has stressed that incorporating intuitive, *human* priors will be necessary for developing robust and interpretable models. Ilyas et al. [13] discuss how classifiers exploit non-intuitive features, and conclude that it is unreasonable to expect classifiers to behave in a way that is meaningful to humans without explicitly encoding such behavior into the model. Still, it remains challenging to encode meaningful, human priors like “have smoother attribution maps” or “treat this group of features similarly” by penalizing the input gradients or the parameters of a model.

In this work, we propose a new method for encoding abstract priors, called *attribution priors*, in which we directly regularize feature attributions of a model during training. We demonstrate that we can encode meaningful domain knowledge as differentiable functions of feature attributions. Furthermore, we introduce a novel feature attribution method - *expected gradients* - which is an extension of integrated gradients [37] that is naturally suited to being regularized under an attribution prior and avoids hyper-parameter choices required by previous methods. Using attribution priors, we build improved deep models for three different prediction tasks. On images, we use our framework to train a deep model that is more interpretable and generalizes better to noisy data by encouraging the model to have piecewise smooth attribution maps over pixels. On gene expression data, we show how to both reduce prediction error and better capture biological signal by encouraging similarity among gene expression features using a graph prior. Finally, on a patient mortality prediction task, we develop a sparser model and improve performance when learning from limited training data by encouraging a skewed distribution of the feature attributions.

2 Attribution Priors

In this section, we formally define the notion of an attribution prior, and give three different example priors for different data types. Let $X \in \mathbb{R}^{n \times p}$ denote a training dataset with labels $y \in \mathbb{R}^o$, where n is the number of samples, p is the number of features, and o is the number of outputs. For some model parameters θ , let $\Phi(\theta, X)$ be a feature attribution method, which is a function of θ and the data X . Let ϕ_i^ℓ be the feature importance of feature i in sample ℓ . In standard deep learning we aim to find optimal parameters θ by minimizing loss, subject to a regularization term $\Omega'(\theta)$ on the parameters:

$$\theta = \operatorname{argmin}_\theta \mathcal{L}(\theta; X, y) + \lambda' \Omega'(\theta).$$

Given a feature attribution method $\Phi(\theta, X)$, we formally define an *attribution prior* as a scalar-valued penalty function of the feature attributions, $\Omega(\Phi(\theta, X))$, which represents a prior probability distribution over possible attributions (see 2.1-2.3 and Supplement for details). When learning our model, we optimize with respect to both training loss and the attribution prior:

$$\theta = \operatorname{argmin}_\theta \mathcal{L}(\theta; X, y) + \lambda \Omega(\Phi(\theta, X)),$$

where λ is the regularization strength. We can define different attribution priors for different tasks depending on the data and our domain knowledge. The attribution prior function Ω is agnostic to the attribution method Φ . While in Section 3 we propose a feature attribution method for attribution priors, other attribution methods can also be used. In particular, by viewing input gradients as explanations, attribution priors are a generalization of existing gradient-based penalties [16, 26, 40, 14, 27].

To demonstrate how attribution priors can capture human intuition in a variety of domains, below we define and apply three example priors for three different data types.

2.1 Pixel Attribution Prior for Image Classification

Prior work on interpreting image models has focused on creating pixel attribution maps, which assign a numerical value to each pixel indicating how important that pixel was for a model’s prediction [31, 37]. These attribution maps can be noisy and often highlight seemingly unimportant pixels in the background. Such attributions can be difficult to understand, and may indicate that an image model is vulnerable to adversarial attacks [25]. Although we may desire a model with smoother attributions, existing methods only post-process attribution maps and do not change model behavior [34, 31, 9]. Such techniques may not be faithful to the original model [13]. In this section, we describe how to apply our framework to train image models with naturally smoother attributions.

To regularize pixel-level attributions, we use the following intuition: neighboring pixels should have a similar impact on an image model’s output. To encode this intuition, we apply a total variation loss on pixel-level attributions as follows [4]:

$$\Omega_{\text{pixel}}(\Phi(\theta, X)) = \sum_{\ell} \sum_{i,j} |\phi_{i+1,j}^{\ell} - \phi_{i,j}^{\ell}| + |\phi_{i,j+1}^{\ell} - \phi_{i,j}^{\ell}|,$$

where $\phi_{i,j}^{\ell}$ is the attribution for the i, j -th pixel in the ℓ -th training image. Including the λ scale factor, this penalty is equivalent to placing a Laplace($0, \lambda^{-1}$) prior on the differences between adjacent pixel attributions. For further details, see [1] and the Supplement.

2.2 Graph Attribution Prior for Gene Expression Data

In the image domain, our attribution prior took the form of a penalty encouraging smoothness over adjacent pixels. In other domains, we may have prior information about specific relationships between features that can be encoded as an arbitrary graph (such as social networks, knowledge graphs, or protein-protein interactions). For example, prior work in bioinformatics has shown that protein-protein interaction networks contain valuable information that can be used to improve performance on biological prediction tasks [5]. These networks can be represented as a weighted, undirected graph. Formally, say we have a weighted adjacency matrix $W \in \mathbb{R}_+^{p \times p}$ for an undirected graph, where the entries encode our prior belief about the pairwise similarity of the importances between two features. For a biological network, $W_{i,j}$ encodes either the probability or strength of interaction between the i -th and j -th genes (or proteins). We can encourage similarity along graph edges by penalizing the squared Euclidean distance between each pair of feature attributions in proportion to how similar we believe them to be. Using the graph Laplacian ($L_G = D - W$), where D is the diagonal degree matrix of the weighted graph this becomes:

$$\Omega_{\text{graph}}(\Phi(\theta, X)) = \sum_{i,j} W_{i,j} (\bar{\phi}_i - \bar{\phi}_j)^2 = \bar{\phi}^T L_G \bar{\phi}.$$

In this case, we choose to penalize *global* rather than local feature attributions. So we define $\bar{\phi}_i$ to be the importance of feature i across all samples in our data set, where this global attribution is calculated as the average magnitude of the feature attribution across all samples: $\bar{\phi}_i = \frac{1}{n} \sum_{\ell=1}^n |\phi_i^{\ell}|$. Overall, Ω_{graph} is equivalent to placing a Normal($0, \lambda^{-1}$) prior on the differences between attributions for features that are adjacent in the graph. We refer the reader to [1] and the Supplement for details.

2.3 Sparsity Attribution Prior for Feature Selection

Feature selection and *sparsity* are popular ways to alleviate the curse of dimensionality, facilitate model interpretability, and improve generalization performance by building models that rely only on a small number of input features. One of the most straightforward ways to build a sparse deep model is to apply an L1 penalty to the first layer (and possibly subsequent layers) of the network. This can be extended in the sparse group lasso (SGL) to penalize all weights connected to a given feature [8, 30]. Another method, as in Ross et al. [24], is to penalize the model’s input gradients.

These approaches suffer from two problems: First, a feature with small first-layer weights or input gradients may still strongly affect the model’s output [33]. A feature whose attribution value (e.g., integrated or expected gradient) is zero, on the other hand, is much less likely to have any affect on predictions. Second, successfully minimizing the L1 or SGL penalty is not necessarily the best way to create a sparse model. A model that puts weight w on 1 feature is penalized more than one that puts weight $\frac{w}{2^p}$ on each of p features. Prior work on sparse linear regression has shown that the Gini coefficient G of the weights, defined as 0.5 minus the area under the CDF of sorted values, avoids such problems and corresponds more directly to a sparse model [12, 42]. We extend this analysis to deep models by noting that the Gini coefficient can be written differentiably. Thus we can use it to develop an attribution penalty based on the global feature attributions $\bar{\phi}_i$:

$$\Omega_{\text{sparse}}(\Phi(\theta, X)) = -2G(\Phi) = -\frac{\sum_{i=1}^p \sum_{j=1}^p |\bar{\phi}_i - \bar{\phi}_j|}{n \sum_{i=1}^p \bar{\phi}_i}.$$

This is similar to the total variation penalty Ω_{image} , but normalized and with a flipped sign to *encourage* differences. The corresponding attribution prior is maximized when global attributions are zero for all but one feature, and minimized when attributions are uniform across features.

3 Expected Gradients

Here we propose a feature attribution method, called *expected gradients*, and describe why it is a natural choice for attribution priors. Expected gradients is an extension of integrated gradients [37] with fewer hyperparameter choices. Like several other attribution methods, integrated gradients aims to explain the difference between a model’s current prediction and the prediction that model would make when given a baseline input. This baseline input is meant to represent some uninformative reference input; that is, it represents not knowing the value of the input features. Although the choice is necessary for several feature attribution methods [37, 33, 3], the choice is often made arbitrarily. For example, in image tasks, the image of all zeros is often chosen as a baseline, but doing so implies that black pixels will not be highlighted as important by existing feature attribution methods. In many domains, it is not clear how to choose a baseline that correctly represents a lack of information.

Our method avoids an arbitrary choice of baseline by modeling not knowing the value of a feature by integrating over a dataset. For a model f , the *integrated gradients* value for feature i is defined as:

$$\text{IntegratedGradients}_i(x, x') := (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\delta f(x' + \alpha \times (x - x'))}{\delta x_i} \delta \alpha,$$

where x is the target input and x' is baseline input. To avoid specifying x' , we define the *expected gradients* value for feature i as:

$$\text{ExpectedGradients}_i(x) := \int_{x'} \left((x_i - x'_i) \int_{\alpha=0}^1 \frac{\delta f(x' + \alpha \times (x - x'))}{\delta x_i} \delta \alpha \right) p_D(x') \delta x',$$

where D is the underlying data distribution. Since expected gradients is also a diagonal path method, it satisfies the same axioms as integrated gradients [10]. Directly integrating over the training distribution is intractable; so we instead reformulate the integrals as expectations:

$$\text{ExpectedGradients}_i(x) := \mathbb{E}_{x' \sim D, \alpha \sim U(0,1)} \left[(x_i - x'_i) \frac{\delta f(x' + \alpha \times (x - x'))}{\delta x_i} \right].$$

This expectation-based formulation lends itself to a natural sampling based approximation method: draw samples of x' from the training dataset and α from $U(0, 1)$, compute the value inside the expectation for each sample, and average over samples. In our datasets, we can get reliable attributions using around 200 samples (see Supplement), similar to integrated gradients [37].

Training with expected gradients: If we let the attribution function Φ in our attribution prior $\Omega(\Phi(\theta, X))$ be expected gradients, this would appear computationally expensive. Computing a good approximation during training may require computing potentially hundreds of extra gradient calls every training step, as discussed in the previous paragraph. Ordinarily, this would make training with such attributions intractable. However, we observe that most deep learning models today are

Table 1: Benchmark results on synthetic data with correlated features. Larger numbers are better for all metrics. For metric names (K = Keep, R = Remove), (P = Positive, N = Negative, A = Absolute), (M = Mean masking, R = Resample masking, and I = Impute masking) (see Supplement for details).

Method	KPM	KPR	KPI	KNM	KNR	KNI	KAM	KAR	KAI
Expected Grad.	3.731	3.800	3.973	3.615	3.551	3.873	0.906	0.903	0.919
Integrated Grad.	3.667	3.736	3.920	3.543	3.476	3.808	0.905	0.899	0.920
Inputs * Grad.	3.633	3.698	3.881	3.473	3.412	3.732	0.904	0.897	0.920
Gradients	0.096	0.122	0.099	0.076	-0.112	0.052	0.838	0.823	0.887
Random	0.033	0.106	0.077	-0.012	-0.093	-0.053	0.593	0.583	0.715

Method	RPM	RPR	RPI	RNM	RNR	RNI	RAM	RAR	RAI
Expected Grad.	3.612	3.575	3.525	3.759	3.830	3.683	0.897	0.885	0.880
Integrated Grad.	3.539	3.503	3.365	3.687	3.754	3.543	0.872	0.859	0.822
Inputs * Grad.	3.474	3.442	3.273	3.655	3.718	3.457	0.869	0.854	0.812
Gradients	0.035	-0.098	-0.020	0.110	0.105	0.108	0.729	0.712	0.616
Random	-0.053	-0.100	-0.106	0.034	0.092	0.111	0.400	0.400	0.275

trained using some variant of batch gradient descent, in which the gradient of a loss function is approximated over many training steps using mini-batches of data. We can use a batch training procedure to approximate expected gradients over the training procedure as well. During training, we let k be the number of samples we draw to compute expected gradients for each mini-batch of data. Remarkably, we find that as small as $k = 1$ suffices to regularize the explanations, which corresponds to only one additional gradient call per training step. More details about how we train with respect to the attribution term can be found in the Supplement.

4 Experiments

We first evaluate expected gradients by comparing it with other feature attribution methods on 18 benchmarks introduced in Lundberg et al. [21] (Table 1). These benchmark metrics aim to evaluate how well each attribution method finds the most important features for a given dataset and model. For all metrics, a larger number corresponds to a better feature attribution method. Expected gradients significantly outperforms the next best feature attribution method ($p = 7.2 \times 10^{-5}$, one-tailed Binomial test). We provide more details on the metrics, and also additional benchmark results for a second dataset in the Supplement.

4.1 A Pixel Attribution Prior Improves Robustness to Image Noise

We apply our Ω_{pixel} attribution prior to the MNIST dataset [16]. We train a CNN model with two convolutional layers and a single hidden layer. We search over λ values logarithmically spaced in the range $[10^{-20}, 10^{-8}]$, and plot both test accuracy and total variation of attributions as a function of λ in Figure 1. The figure demonstrates that we can achieve equivalent test accuracy (99.55% for the baseline model vs. 99.52% for the regularized model) while making attribution maps on test images smoother. The figure also depicts the attribution maps for the baseline and the regularized model at $\lambda = 5 \times 10^{-11}$. The regularized attribution maps more clearly highlight the digit.

Recent work in understanding image classifiers has suggested that they are brittle to small domain shifts: small changes in the underlying distribution of the training and test set, such as from data cleaning procedures, can result in significant drops in test accuracy [23]. To simulate such a domain shift, we randomly set a fraction of pixels in the test images to be 0 or 1 with equal probability and then re-evaluated both our regularized model and the baseline model. In Figure 1, we plot both the result and an example image at each noise level. The plot shows that our regularized model better generalizes to a shift in the underlying data distribution at every noise level. In particular, at 10% pixel noise, a level which is still easily read by humans, the baseline model drops below 50% accuracy, while our regularized model maintains an accuracy above 90%.

For more details on our training procedures and architectures, as well as more examples of attribution maps and experiments on the ImageNet dataset [28], please refer to the Supplement.

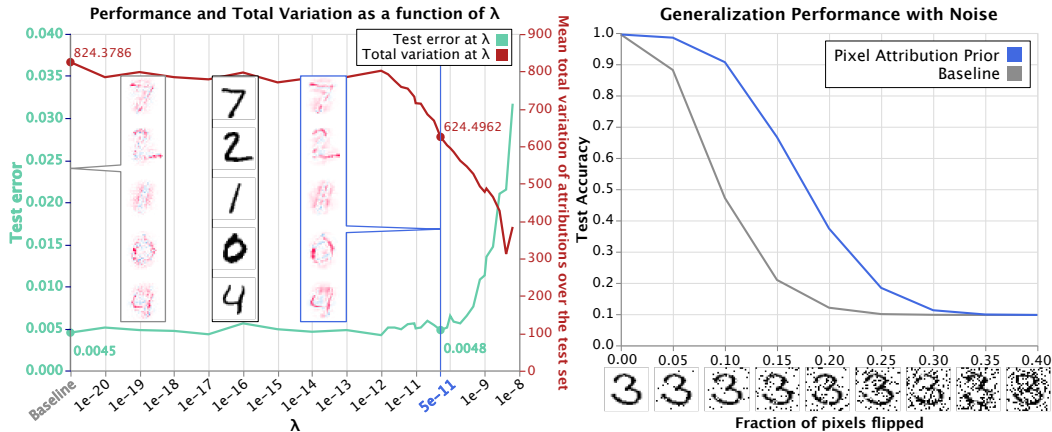


Figure 1: Penalizing differences between feature attributions of adjacent pixels creates more understandable image explanations. Left: A plot of attribution variation and test error vs. regularization strength, as well as feature attributions for the baseline model (gray call-out box) and the regularized model at $\lambda = 5 \times 10^{-11}$ (blue call-out box) for the true output class. Red pixels contribute positively to the true class, while blue pixels contribute negatively. Right: A plot of our model at $\lambda = 5 \times 10^{-11}$ vs. the baseline model predicting on progressively noisier test images. Models with smooth attributions generalize better to noise.

4.2 A Graph Attribution Prior Improves Anti-Cancer Drug Response Prediction

In this section, we show how incorporating the Ω_{graph} attribution prior not only leads to a model with more reasonable attributions, but also improves predictive performance by allowing us to incorporate prior biological knowledge into the training process. We downloaded publicly available gene expression and drug response data for patients with acute myeloid leukemia (AML, a type of blood cancer) and tried to predict patients’ drug response from their gene expression [38]. For this regression task, an input sample was a patient’s gene expression profile plus a one-hot encoded vector indicating which drug was tested in that patient, while the label we tried to predict was drug response (measured by IC50 - the concentration of the drug required to kill half of the patient’s tumor cells). To define the graph used by our prior we downloaded the tissue-specific gene interaction graph for the tissue most closely related to AML in the HumanBase database [11].

We find that a two-layer neural network trained with our graph attribution prior (Ω_{graph}) significantly outperforms all other methods in terms of test set performance as measured by R^2 (Figure 2). Unsurprisingly, when we replace the biological graph from HumanBase with a randomized graph, we find that the test performance is no better than the performance of a neural network trained without *any* attribution prior. We also observe significantly improved test performance when using the prior graph information to regularize a linear LASSO model. Finally, we note that our graph attribution prior neural network significantly outperforms a recent method for utilizing graph information in deep neural networks, graph convolutional neural networks [15].

To see if our model’s attributions match biological intuition we conducted Gene Set Enrichment Analysis (a modified Kolmogorov–Smirnov test) to see if our top genes, as ranked by mean absolute feature attribution, were enriched for membership in any pathways (see the Supplement for more details, including the top pathways for each model) [36]. We see that the neural network with the tissue-specific graph attribution prior captures significantly more biologically-relevant pathways (increased number of significant pathways after FDR correction using the Benjamini-Hochberg procedure) than a neural network without attribution priors (See Figure 2) [2]. Furthermore, the pathways used by the model with the graph attribution prior more closely match with biological expert knowledge – pathways included prognostically useful AML gene expression profiles, as well as important AML-related transcription factors (see Figure 2 and Supplement) [17, 39].

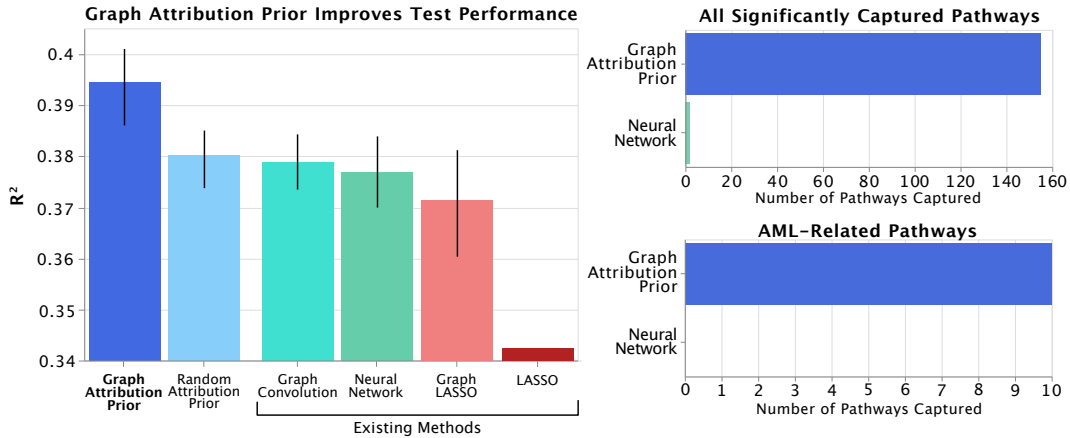


Figure 2: Left: A neural network trained with our graph attribution prior attains the best test performance. Top Right: A neural network trained with our graph attribution prior has far more significantly captured biological pathways than a standard neural network. Bottom Right: The graph attribution prior network captures more pathways, and also captures more AML-relevant pathways.

4.3 A Sparsity Prior Improves Performance with Limited Training Data

Here, we show that the Ω_{sparse} attribution prior can build sparser models that perform significantly better in settings with limited training data. We use a publicly available healthcare mortality prediction dataset of 13,000 patients [22], where each of the 36 features (119 after one-hot encoding) represent medical data such as a patient’s age, vital signs, and laboratory measurements. The outcome is a binary label indicating survival after 10 years. Sparse models in this setting may reduce cost by accurately risk-stratifying patients using fewer lab tests, or may enable accurate models to be trained with very few labeled patient samples. We build 3-layer binary classifier neural networks regularized using L1, sparse group lasso (SGL) and sparse attribution prior penalties to predict patient survival. The regularization strength was determined by a broad sweep (from 10^{-10} to 10^3 for all methods) followed by a fine sweep in the area of lowest loss (see Supplement for details).

We plot the sparsity of the model (Gini coefficient) against its test ROC-AUC score for a range of regularization strengths (Figure 3). While all models can achieve an ROC-AUC of 0.85, the sparse attribution prior attains substantially higher sparsity at this performance point (Figure 3 top left). We also plot the cumulative importance of features (sorted in increasing order) for the best performing models – those closest to the top right of the scatterplot (ROC-AUC=1, Gini=1). The resulting curves show that the sparse attribution prior is much more effective at concentrating importance in the top few features (Figure 3 top right). We also show this sparsity enables better predictions when little training data is available by training on a subsample of the data with $n = 100$ for train and validation sets (Figure 3 bottom). To address randomness in subsampling and random starting, we train 100 models of each type on 34 random subsamples of the data and plot average test performance of the top 10 models on each subsample. The sparse attribution prior attains higher ROC-AUC scores than unregularized ($p = 1.66 \times 10^{-3}$, $T = 4.25$), L1-regularized ($p = 2.60 \times 10^{-3}$, $T = 4.09$), and SGL-regularized ($p = 0.0266$, $T = 2.32$) models by paired-samples T -test on all replicates.

5 Related Work

There have been many previous attribution methods proposed for deep learning models [18, 3, 33, 37]. We specifically chose to extend integrated gradients because it was easily differentiable and comes with axiomatic justification [37].

Training with gradient penalties has also been discussed by existing literature. Drucker and Le Cun [7] introduced the idea of regularizing the magnitude of input gradients in order to improve generalization performance on digit classification. Since then, gradient regularization has been used extensively as an adversarial defense mechanism in order to minimize changes to network outputs over small perturbations of the input [14, 40, 27]. Ross and Doshi-Velez [25] make a connection between

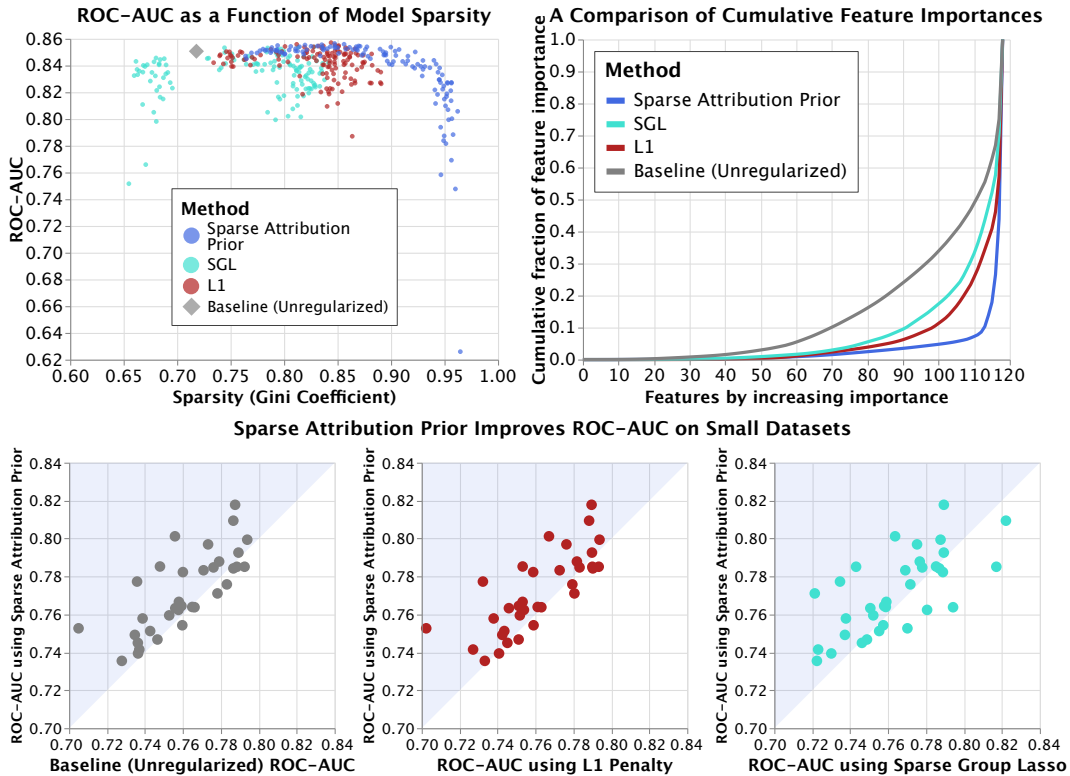


Figure 3: A sparsity-inducing feature attribution prior leads to sparse models with better performance in small-data settings. Top left: A sparse attribution prior increases sparsity while preserving ROC-AUC. Each dot represents a model with a different regularization strength. Top right: A sparse attribution prior concentrates a larger fraction of global feature importance in the top few features. Bottom: A sparse attribution prior enables more accurate predictions across 34 small subsampled datasets (100 samples each). Each dot is performance on a different subsampled dataset.

gradient-based training for adversarial purposes and network interpretability. Ilyas et al. [13] formally describe how the phenomena of adversarial examples may arise due to features that are predictive yet non-intuitive, and stress the need to incorporate human intuition into the training process if we expect models to behave according to such intuition.

There is very little previous work on actually incorporating feature attribution methods into training. Sen et al. [32] formally describe the problem of classifiers having unexpected behavior on inputs not seen in the training distribution, like those generated by asking whether a prediction would change if a particular feature value changed. They describe an active learning algorithm that updates a model based on points generated from a counter-factual distribution. Their work differs from ours in that they use feature attributions to generate counter-factual examples, but do not directly penalize the attributions themselves. Ross et al. [26] introduce the idea of training models to have correct explanations, not just good performance. Their method differs from ours in that they use input gradients rather than more modern axiomatically based attribution methods. Their method also differs from ours in that their regularization penalty is binary: it is limited to only encouraging individual features to either be or not be important, and does not consider interactions between attributions.

In the context of image attributions, there has been a line of work in generating sharper and more interpretable attribution maps. Smilkov et al. [34] develop a method for generating sharper attribution maps for image classification networks that works on top of methods like integrated gradients. Selvaraju et al. [31] propose a method based on back-propagating gradients from an output logit corresponding to the class of interest that achieves sharp, localized visualizations. Fong and Vedaldi [9] design an algorithm to learn a map of which pixels, when perturbed, most change the output of a classifier. All of these methods generate smooth attribution maps. However, they do not change the

training procedure, meaning the underlying model being explained may still suffer from problems like placing too much importance on individual pixels.

6 Discussion

The immense popularity of deep learning has driven its application in a wide range of domains. Each of these domains, be it image analysis, biology, or healthcare, involves different types of potentially complicated prior knowledge. While it is in principle possible to hand-design network architectures to encode this prior knowledge, we propose a simpler approach. Using attribution priors, any knowledge that can be encoded as a differentiable function of feature attributions can be used to encourage a model to act in a particular way in a particular domain.

By introducing expected gradients, we demonstrate an axiomatically based feature attribution method that enables the natural integration of attribution priors into stochastic gradient descent. Expected gradients also removes the choice of a single reference value that many existing feature attribution methods require, instead using the training distribution as a reference.

Training with expected gradients allows us to improve model performance by encoding prior knowledge across several different domains. This leads to smoother and more interpretable image models, biological predictive models that incorporate graph-based prior knowledge, and sparser health care models that can perform better in data-scarce scenarios. Attribution priors provide a broadly applicable framework for encoding domain knowledge, and we believe they will be valuable across a wide array of domains in the future.

Acknowledgments

This work was funded by National Science Foundation [DBI-1759487, DBI-1552309]; American Cancer Society [127332-RSG-15-097-01-TBG]; and National Institutes of Health [R35 GM 128638, R01 AG 061132].

References

- [1] J. M. Bardsley. Laplace-distributed increments, the laplace prior, and edge-preserving regularization. *J. Inverse Ill-Posed Probl*, 2012.
- [2] Y. Benjamini and Y. Hochberg. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300, 1995. ISSN 00359246. URL <http://www.jstor.org/stable/2346101>.
- [3] A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, and W. Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *International Conference on Artificial Neural Networks*, pages 63–71. Springer, 2016.
- [4] A. Chambolle. An algorithm for total variation minimization and applications. *Journal of Mathematical imaging and vision*, 20(1-2):89–97, 2004.
- [5] W. Cheng, X. Zhang, Z. Guo, Y. Shi, and W. Wang. Graph-regularized dual Lasso for robust eQTL mapping. *Bioinformatics*, 30(12):i139–i148, 06 2014. ISSN 1367-4803. doi: 10.1093/bioinformatics/btu293. URL <https://doi.org/10.1093/bioinformatics/btu293>.
- [6] A. Datta, S. Sen, and Y. Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In *Security and Privacy (SP), 2016 IEEE Symposium on*, pages 598–617. IEEE, 2016.
- [7] H. Drucker and Y. Le Cun. Improving generalization performance using double backpropagation. *IEEE Transactions on Neural Networks*, 3(6):991–997, 1992.
- [8] J. Feng and N. Simon. Sparse-input neural networks for high-dimensional nonparametric regression and classification. *arXiv preprint arXiv:1711.07592*, 2017.
- [9] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3429–3437, 2017.
- [10] E. J. Friedman. Paths and consistency in additive cost sharing. *International Journal of Game Theory*, 32(4):501–518, 2004.
- [11] C. S. Greene, A. Krishnan, A. K. Wong, E. Ricciotti, R. A. Zelaya, D. S. Himmelstein, R. Zhang, B. M. Hartmann, E. Zaslavsky, S. C. Sealfon, et al. Understanding multicellular function and disease with human tissue-specific networks. *Nature genetics*, 47(6):569, 2015.
- [12] N. Hurley and S. Rickard. Comparing measures of sparsity. *IEEE Transactions on Information Theory*, 55(10):4723–4741, 2009.
- [13] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019.
- [14] D. Jakubovitz and R. Giryes. Improving dnn robustness to adversarial attacks using jacobian regularization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 514–529, 2018.
- [15] T. N. Kipf and M. Welling. Semi-Supervised Classification with Graph Convolutional Networks. *CoRR*, abs/1609.0, 2016. URL <http://arxiv.org/abs/1609.02907>.
- [16] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2:18, 2010.
- [17] J. Liu, Y.-Z. Qin, S. Yang, Y. Wang, Y.-J. Chang, T. Zhao, Q. Jiang, and X.-J. Huang. Meis1 is critical to the maintenance of human acute myeloid leukemia cells independent of MLL rearrangements. *Annals of Hematology*, 96(4):567–574, apr 2017. ISSN 1432-0584. doi: 10.1007/s00277-016-2913-6. URL <https://doi.org/10.1007/s00277-016-2913-6>.
- [18] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.

- [19] S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent individualized feature attribution for tree ensembles. arXiv preprint arXiv:1802.03888, 2018.
- [20] S. M. Lundberg, B. Nair, M. S. Vavilala, M. Horibe, M. J. Eisses, T. Adams, D. E. Liston, D. K.-W. Low, S.-F. Newman, J. Kim, et al. Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. Nature Biomedical Engineering, 2(10):749, 2018.
- [21] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmel-farb, N. Bansal, and S.-I. Lee. Explainable ai for trees: From local explanations to global understanding, 2019.
- [22] H. W. Miller. Plan and operation of the health and nutrition examination survey, united states, 1971-1973. DHEW publication no.(PHS)-Dept. of Health, Education, and Welfare (USA), 1973.
- [23] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar. Do imagenet classifiers generalize to imagenet? arXiv preprint arXiv:1902.10811, 2019.
- [24] A. Ross, I. Lage, and F. Doshi-Velez. The neural lasso: Local linear sparsity for interpretable explanations. In Workshop on Transparent and Interpretable Machine Learning in Safety Critical Environments, 31st Conference on Neural Information Processing Systems, 2017.
- [25] A. S. Ross and F. Doshi-Velez. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In Thirty-second AAAI conference on artificial intelligence, 2018.
- [26] A. S. Ross, M. C. Hughes, and F. Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. arXiv preprint arXiv:1703.03717, 2017.
- [27] K. Roth, A. Lucchi, S. Nowozin, and T. Hofmann. Adversarially robust training through structured gradient regularization. arXiv preprint arXiv:1805.08736, 2018.
- [28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. International journal of computer vision, 115(3):211–252, 2015.
- [29] R. Sayres, A. Taly, E. Rahimy, K. Blumer, D. Coz, N. Hammel, J. Krause, A. Narayanaswamy, Z. Rastegar, D. Wu, et al. Using a deep learning algorithm and integrated gradients explanation to assist grading for diabetic retinopathy. Ophthalmology, 126(4):552–564, 2019.
- [30] S. Scardapane, D. Comminiello, A. Hussain, and A. Uncini. Group sparse regularization for deep neural networks. Neurocomputing, 241:81–89, 2017.
- [31] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE International Conference on Computer Vision, pages 618–626, 2017.
- [32] S. Sen, P. Mardziel, A. Datta, and M. Fredrikson. Supervising feature influence. arXiv preprint arXiv:1803.10815, 2018.
- [33] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating activation differences. In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 3145–3153. JMLR. org, 2017.
- [34] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg. Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825, 2017.
- [35] E. Štrumbelj and I. Kononenko. Explaining prediction models and individual predictions with feature contributions. Knowledge and information systems, 41(3):647–665, 2014.
- [36] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, and J. P. Mesirov. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. Proceedings of the National Academy of Sciences, 102(43):15545 LP – 15550, oct 2005. doi: 10.1073/pnas.0506580102. URL <http://www.pnas.org/content/102/43/15545.abstract>.

- [37] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 3319–3328. JMLR. org, 2017.
- [38] J. W. Tyner, C. E. Tognon, D. Bottomly, B. Wilmot, S. E. Kurtz, S. L. Savage, N. Long, A. R. Schultz, E. Traer, M. Abel, et al. Functional genomic landscape of acute myeloid leukaemia. Nature, 562(7728):526, 2018.
- [39] P. J. M. Valk, R. G. W. Verhaak, M. A. Beijen, C. A. J. Erpelinck, S. B. v. W. van Doorn-Khosrovani, J. M. Boer, H. B. Beverloo, M. J. Moorhouse, P. J. van der Spek, B. Löwenberg, and R. Delwel. Prognostically Useful Gene-Expression Profiles in Acute Myeloid Leukemia. New England Journal of Medicine, 350(16):1617–1628, 2004. doi: 10.1056/NEJMoa040465. URL <https://doi.org/10.1056/NEJMoa040465>.
- [40] F. Yu, Z. Xu, Y. Wang, C. Liu, and X. Chen. Towards robust training of neural networks by regularizing adversarial gradients. arXiv preprint arXiv:1805.09370, 2018.
- [41] J. R. Zech, M. A. Badgeley, M. Liu, A. B. Costa, J. J. Titano, and E. K. Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study. PLoS medicine, 15(11):e1002683, 2018.
- [42] D. Zonoobi, A. A. Kassim, and Y. V. Venkatesh. Gini index as sparsity measure for signal reconstruction from compressive samples. IEEE Journal of Selected Topics in Signal Processing, 5(5):927–932, 2011.

Learning Explainable Models using Attribution Priors - Supplement

1 Training with Attributions

Normally, training with a penalty on any function of the gradients would require solving a differential equation. To avoid this, we adopt a double back-propagation scheme in which the gradients are first calculated with respect to the training loss, and alternately calculated with the loss with respect to the attributions [25, 3]. In Tensorflow, we accomplish this with two different training operations: one to minimize the training loss, and one to minimize the loss with respect to the attributions [1].

Our attribution method, expected gradients, requires background reference samples to be drawn from the training data. More specifically, for each input in a batch of inputs, we need k additional inputs to calculate expected gradients for that input batch. As long as k is smaller than the batch size, we can avoid any additional data reading by re-using the same batch of input data as a reference batch, as in Zhang et al. [26]. We accomplish this by shifting the batch of input k times, such that each input in the batch uses k other inputs from the batch as its reference values.

2 Corresponding priors

In this section, we elaborate on the explicit form of the attribution priors used in the paper. In general, minimizing the error of a model corresponds to maximizing the likelihood of the data under a generative model consisting of the learned model plus parametric noise. For example, minimizing mean squared error in a regression task corresponds to maximizing the likelihood of the data under the learned model, assuming Gaussian-distributed errors.

$$\operatorname{argmin}_{\theta} \|f_{\theta}(X) - y\|_2^2 = \operatorname{argmax}_{\theta} \exp(-\|f_{\theta}(X) - y\|_2^2) = \theta_{MLE}$$

where θ_{MLE} is the maximum-likelihood estimate of θ under the model $Y = f_{\theta}(X) + \mathcal{N}(0, \sigma)$.

An additive regularization term is equivalent to adding a multiplicative (independent) prior to yield a maximum a posteriori estimate:

$$\operatorname{argmin}_{\theta} \|f_{\theta}(X) - y\|_2^2 + \lambda \|\theta\|_2^2 = \operatorname{argmax}_{\theta} \exp(-\|f_{\theta}(X) - y\|_2^2) \exp(-\lambda \|\theta\|_2^2) = \theta_{MAP}$$

Here adding an L2 penalty is equivalent to MAP for $Y = f_{\theta}(X) + \mathcal{N}(0, \sigma)$ with a $\mathcal{N}(0, \frac{1}{\lambda})$ prior. The natural next question is what attribution priors are being enforced by the penalties used in our experiments.

Image prior: Our image prior uses a total variation penalty, which has been well-studied.

$$\Omega_{\text{pixel}}(\Phi(\theta, X)) = \sum_{\ell} \sum_{i,j} |\phi_{i+1,j}^{\ell} - \phi_{i,j}^{\ell}| + |\phi_{i,j+1}^{\ell} - \phi_{i,j}^{\ell}|$$

It has been shown in Bardsley [2] that this penalty is equivalent to placing 0-mean, iid, Laplace-distributed priors on the differences between adjacent pixel values. That is, $\phi_{i+1,j}^{\ell} - \phi_{i,j}^{\ell} \sim \text{Laplace}(0, \lambda^{-1})$ and $\phi_{i,j+1}^{\ell} - \phi_{i,j}^{\ell} \sim \text{Laplace}(0, \lambda^{-1})$. Bardsley [2] does not call our penalty "total variation", but it is in fact the widely used anisotropic version of total variation, and is directly implemented in Tensorflow [1, 9, 18].

Graph prior: The graph prior extends the image prior to arbitrary graphs:

$$\Omega_{\text{graph}}(\Phi(\theta, X)) = \bar{\phi}^T L_G \bar{\phi}$$

Just as the image penalty is equivalent to placing a Laplace prior on adjacent pixels in a regular graph, the graph penalty Ω_{graph} is equivalent to placing a Gaussian prior on adjacent features in an arbitrary graph with Laplacian L_G [2].

Sparsity prior: Our sparsity prior uses the Gini coefficient as a penalty, which is written

$$\Omega_{\text{sparse}}(\Phi(\theta, X)) = -\frac{\sum_{i=1}^p \sum_{j=1}^p |\bar{\phi}_i - \bar{\phi}_j|}{n \sum_{i=1}^p \bar{\phi}_i} = -2G(\Phi)$$

By taking exponentials of this function, we find that minimizing the sparsity regularizer is equivalent to maximizing likelihood under a prior proportional to the following:

$$\prod_{i=1}^p \prod_{j=1}^p \exp\left(\frac{1}{\sum_{i=1}^p \bar{\phi}_i} |\bar{\phi}_i - \bar{\phi}_j|\right)$$

To our knowledge, this prior does not directly correspond to a named distribution. However, we can note that its maximum value occurs when one $\bar{\phi}_i$ is 1 and all others are 0, as well as that its minimum occurs when all $\bar{\phi}_i$ are equal.

3 Benchmarking Expected Gradients

3.1 Sampling convergence

Since expected gradients reformulates feature attribution as an *expected value* over two distributions (where background samples x' are drawn from the data distribution and the linear interpolation parameter α is drawn from $U(0, 1)$), we wanted to ensure that we are drawing an adequate number of background samples for convergence of our attributions when benchmarking the performance of our attribution method. Since our benchmarking code was run on the Correlated Groups 60 synthetic dataset, as a baseline we explain all 1000 samples of this data sampling the full dataset (1000 samples) as background samples. To assess convergence to the attributions attained at this number of samples, we measure the mean absolute difference between the attribution matrices resulting from different numbers of background samples (see Figure 1). We empirically find that our attributions are well-converged by the time 100-200 background samples are drawn. Therefore, for the rest of our benchmarking experiments, we used 200 as the number of background samples.

3.2 Benchmark evaluation metrics

To compare the performance of expected gradients with other feature attribution methods, we used the benchmark metrics proposed in Lundberg et al. [10]. These metrics were selected as they are comprehensive of a variety of recent approaches to evaluating feature importance estimates. For example, the Keep Positive Mask metric (KPM) is used to test how well an attribution method can find the features that lead to the greatest increase in the model’s output. This metric progressively removes features by masking with their mean value, in order from least positive impact on model output to most positive impact on model output, as ranked by the attribution method being evaluated. As more features are masked, the model’s output is decreased, creating a curve. The KPM metric measures the area under this curve (larger area corresponds to better attribution method). In addition to the KPM metric, 17 other similar metrics (e.g. Remove Absolute Resample, Keep Negative Impute, etc.) were used (see supplementary material of Lundberg et al. [10] for more details on benchmark metrics). For all of these metrics, a larger number corresponds to a better attribution method. In addition to finding that Expected Gradients outperforms all other attribution methods on nearly all metrics tested for the dataset shown in Table 1 in the main text (the synthetic Correlated Groups 60 dataset proposed in Lundberg et al. [10]), we also tested all 18 metrics on another dataset proposed in the same paper (Independent Linear 60) and find that Expected Gradients is chosen as the best method by all metrics in that case as well (see Table 1). The Independent Linear 60 dataset is comprised of 60 features, where each feature is a 0-mean, unit variance gaussian random variable plus gaussian noise, and the label to predict is a linear function of these features. The Correlated Groups 60 dataset is essentially the same, but now certain groups of 3 features have 0.99 correlation.

Selecting adequate sample number

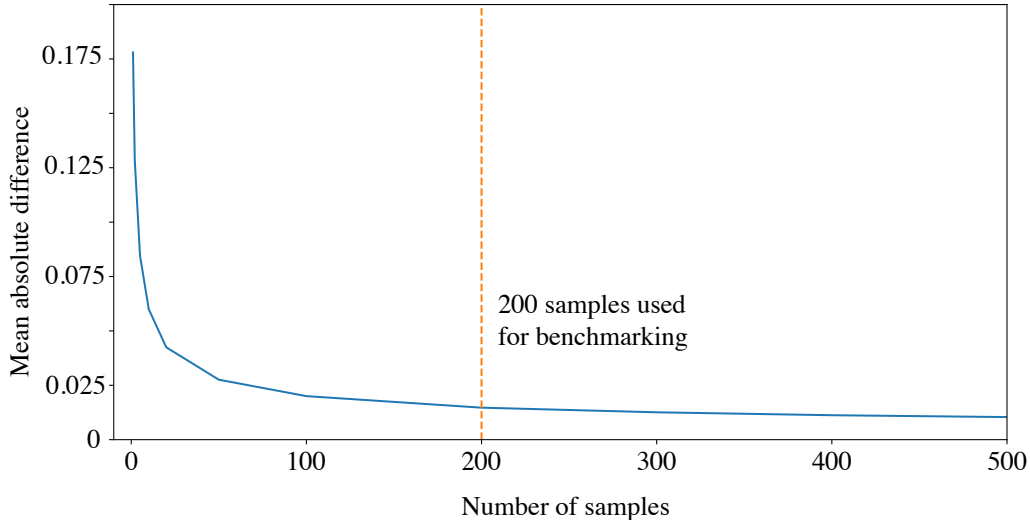


Figure 1: Feature attribution values attained using expected gradients converge as the number of background samples drawn is increased.

Table 1: Benchmark on Independent Linear 60 dataset

Attribution Method	KPM	KPR	KPI	KNM	KNR	KNI	KAM	KAR	KAI
Expected Gradients	4.096	4.179	4.264	4.014	3.835	4.153	0.941	0.946	0.938
Integrated Gradients	4.055	4.112	4.176	3.949	3.753	4.070	0.941	0.945	0.938
Inputs * Gradients	3.992	4.066	4.122	3.896	3.715	4.023	0.939	0.942	0.935
Gradients	0.044	0.107	0.029	0.155	-0.150	0.172	0.902	0.905	0.902
Random	-0.152	0.102	-0.152	0.111	-0.126	0.060	0.470	0.482	0.438
Attribution Method	RPM	RPR	RPI	RNM	RNR	RNI	RAM	RAR	RAI
Expected Gradients	4.079	3.941	4.210	4.203	4.260	4.356	0.992	0.977	1.019
Integrated Gradients	4.013	3.854	4.113	4.157	4.186	4.259	0.973	0.966	0.995
Inputs * Gradients	3.960	3.820	4.065	4.094	4.137	4.199	0.984	0.971	1.009
Gradients	0.110	-0.125	0.133	0.057	0.080	0.041	0.947	0.936	0.985
Random	0.012	-0.124	0.059	0.035	0.101	0.070	0.504	0.521	0.527

For attribution methods to compare, we considered expected gradients (as described in the main text), integrated gradients (as described in Sundararajan et al. [22]), inputs \times gradients (where the attribution for each feature for each sample is the input gradient times the value of that feature in the given sample), gradients, and random.

4 Expected Gradients on ImageNet

One unfortunate consequence of choosing an arbitrary baseline point for methods like integrated gradients is that the baseline point by definition is unimportant. That is, if a user chooses the constant black image as the baseline input, then purely black pixels will not be highlighted as important by integrated gradients. This is true for any constant baseline input. Since expected gradients integrates over a dataset as its baseline input, it avoids forcing a particular pixel value to be unimportant. To demonstrate this, we use the inception v4 network trained on the ImageNet 2012 challenge [23, 16]. We restore pre-trained weights from the Tensorflow Slim library [19]. In Figure 2, we plot attribution

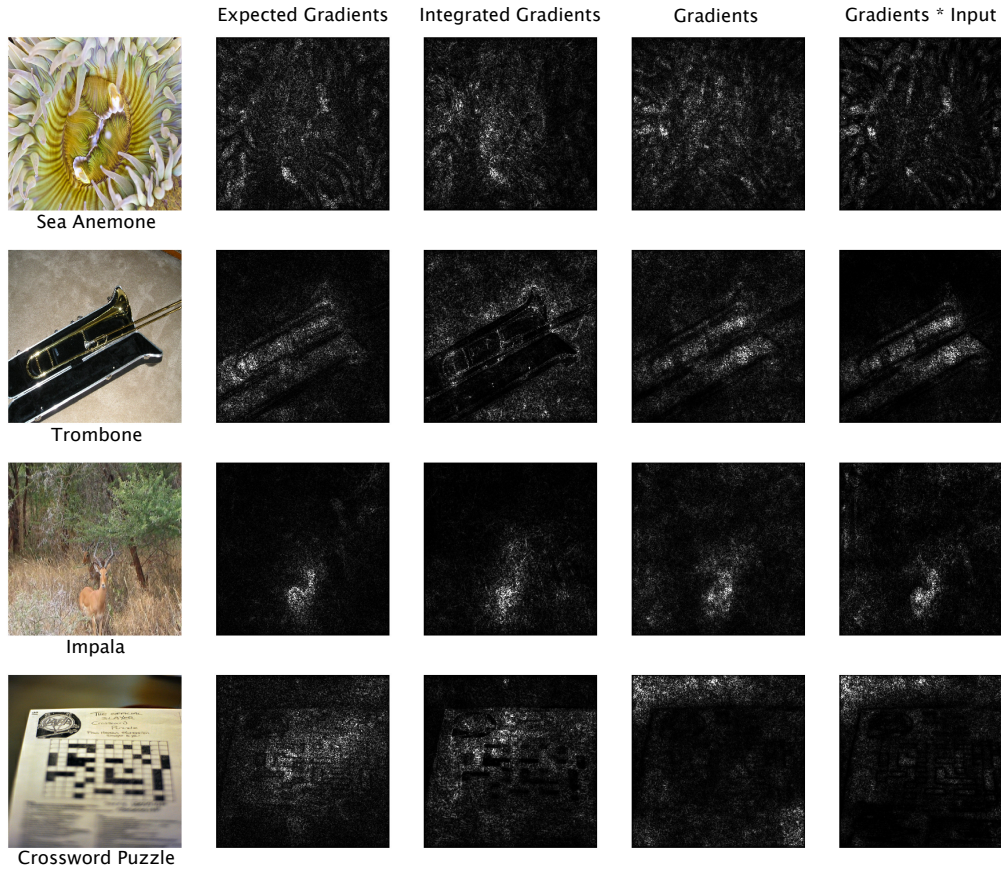


Figure 2: A comparison of attribution methods on ImageNet. Integrated gradients fails to highlight black pixels as important when black is used as a baseline input.

maps of both expected gradients and integrated gradients as well as raw gradients and gradients multiplied by input. Here, we use the constant black image as a baseline input for integrated gradients. For both attribution methods, we use 200 sample/interpolation points. The figure demonstrates that integrated gradients fails to highlight black pixels.

5 MNIST Experiments

5.1 Experimental Setup

On MNIST, we train a CNN with two convolutional layers and a single hidden layer. The convolutional layers have 5x5 filters, a stride length of 1, and 32 and 64 filters total, respectively. Each convolutional layer is followed by a max pooling layer of size 2 with stride length 2. The hidden layer has 1024 units, and a dropout rate of 0.5 during training [21]. Dropout is turned off when calculating the gradients with respect to the attributions. We train with the ADAM optimizer with the default parameters ($\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$) [6]. We train with an initial learning rate of 0.0001, with an exponential decay 0.95 for every epoch, for a total of 60 epochs. We split the MNIST training set randomly into a training set of 45,000 images and 5,000 validation images. Every 500 batches, we evaluate the performance of the network on the validation set. At test time, we restore the model weights that achieved the highest validation performance. For all models, we train with a batch size of 50 images, and use $k = 1$ background reference sample per attribution while training. When training with attributions over images, we first normalize the per-pixel attribution maps by dividing by the standard deviation before computing the total variation - otherwise, the total variation can be made arbitrarily small without changing model predictions by scaling down the pixel attributions close to 0.

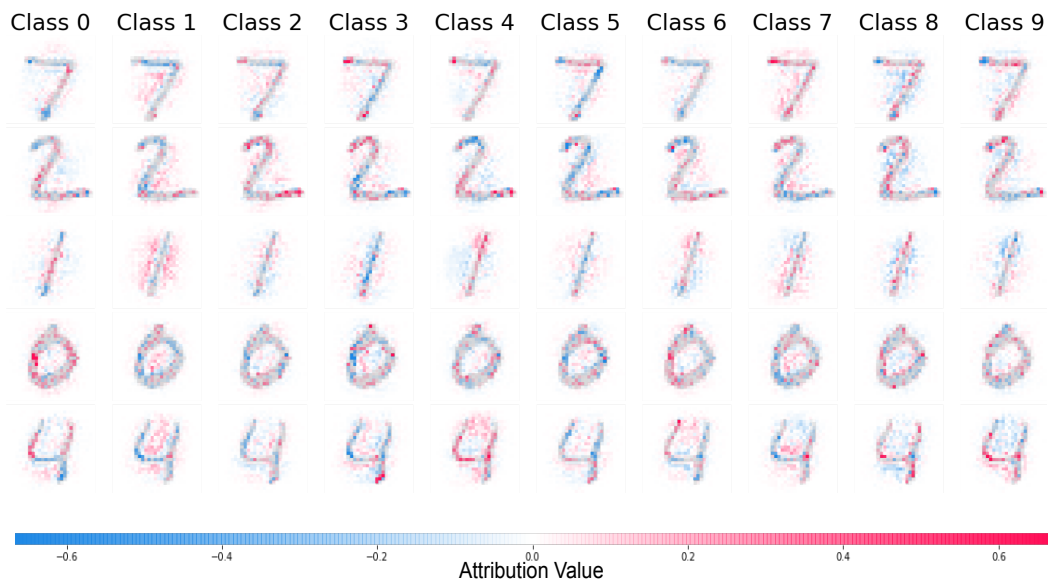


Figure 3: Pixel attribution maps for all classes (classes 0 to 9, left to right) for the baseline model.

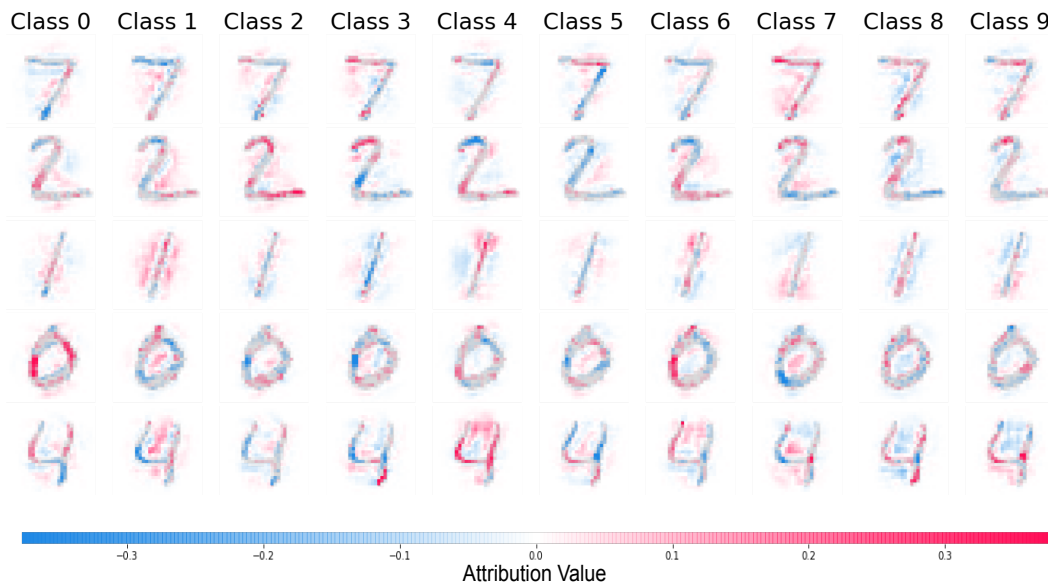


Figure 4: Pixel attribution maps for all classes (classes 0 to 9, left to right) for the regularized model, $\lambda = 5 * 10^{-11}$

5.2 Attribution Maps

In the main text, we plot attribution maps for 5 randomly selected test images for the output class corresponding to the label of the image. In Figure 3, we plot the attribution maps for the same test images for all output classes for the baseline model, and in Figure 4 we do the same for the regularized model. Note that there are patterns that more clearly show up in the regularized attribution maps: for example, the parts of a "2" that contribute to output class "3", and the parts of a "0" that make it appear like a "6".

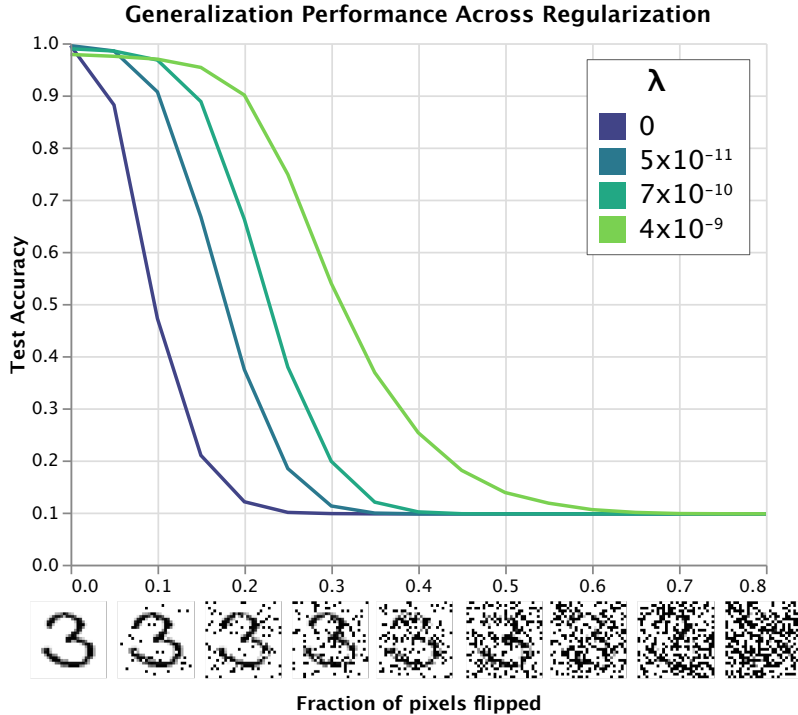


Figure 5: How smooth regularization impacts performance on a noisy test set. The figure shows that as we regularize models to be more smooth (larger λ), they tend to do worse on the original test set, but better when noise is introduced.

5.3 Generalization to Noise

In the main text, we plotted the the generalization performance of the baseline model and the regularized model at $\lambda = 5 \times 10^{-11}$. Here, we include two more values in the curve of test performance vs. noise level, $\lambda = 7 \times 10^{-11}$ and 4×10^{-11} . For clarity, we omit other values of λ . The plot, in Figure 5, demonstrates two trends. First: the more we regularize with λ , the more test performance on the original test set with no noise drops. Second, the same, smooth models that have lower accuracy on the original test set have higher accuracy on the noisy test set. In line with Ilyas et al. [5], this suggests that highly-accurate image classifiers rely not only on features that are intuitive for humans (like the shape and orientation of the lines that make up the digits), but also other, less intuitive and high-frequency features that get corrupted with the introduction of noise.

6 ImageNet Experiments

In this section, we detail experiments performed on applying Ω_{pixel} to classifiers trained on the ImageNet 2012 challenge [16]. We omit this section from the main text since, for computational reasons, the hyper-parameters chosen in this section may not necessarily be optimal.

6.1 Experimental Setup

We use the VGG16 architecture introduced by Simonyan and Zisserman [20]. For computational reasons, we do not train a model from scratch - instead, we fine-tune using pre-trained weights from the Tensorflow Slim package [19]. We fine-tune on the ImageNet 2012 training set using the original cross entropy loss function in addition to Ω_{pixel} using asynchronous gradient updates with a batch size of 16 split across 4 Nvidia 1080 Ti GPUs. During fine-tuning, we use the same training procedure outlined by Silberman and Guadarrama [19]. This includes randomly cropping training images to 224×224 pixels, randomly flipping images horizontally, and normalizing each image to the same range. To optimize, we use gradient descent with a learning rate of 0.00001 and a momentum of 0.9. We use a weight decay of 0.0005, and set $\lambda = 0.00001$ for the first epoch of fine-tuning, and

Table 2: Performance of the VGG16 architecture on the ImageNet 2012 validation dataset before and after fine-tuning.

Model	Top 1 Accuracy	Top 5 Accuracy
Baseline	0.709	0.898
Image Attribution Prior 1 Epoch	0.699	0.886
Image Attribution Prior 1.25 Epochs	0.674	0.876

$\lambda = 0.00002$ for the second epoch of fine-tuning. As with the MNIST experiments, we normalize the feature attributions before taking total variation.

6.2 Results

We plot the attribution maps on images from the validation set using expected gradients for the original VGG16 weights (Baseline), as well as fine-tuned for 320,292 steps (Image Attribution Prior 1 Epoch) and fine-tuned for 382,951 steps, in which the last 60,000 steps were with twice the λ penalty (Image Attribution Prior 1.25 Epochs). Figure 6 demonstrates that fine-tuning using our penalty results in sharper and more interpretable image maps than the baseline network. In addition, we also plot the attribution maps generated by three other methods: integrated gradients (Figure 7), gradients \times inputs (Figure 8) and raw input gradients (Figure 9). Networks regularized with our attribution prior show more clear attribution maps using any of the above methods, which implies that the network is actually viewing pixels more smoothly, independent of the attribution method chosen.

We note that in practice, we observe similar trade-offs between test accuracy and interpretability/robustness mentioned in Ilyas et al. [5]. We show the validation performance of the VGG16 network before and after fine-tuning in Table 2 and observe that the validation accuracy does decrease. However, due to the computational cost even of fine-tuning on ImageNet, we did not perform a hyper-parameter search for the optimal learning rate or λ penalty. We anticipate that with more time and computational resources, we could achieve a better trade-off between interpretable attribution maps and test accuracy.

7 Biological experiments

7.1 RNA-seq preprocessing

To ensure a quality signal for prediction while removing noise and batch effects, it is necessary to carefully preprocess RNA-seq gene expression data. For the biological data experiments, RNA-seq were preprocessed as follows:

1. First, raw transcript counts were converted to fragments per kilobase of exon model per million mapped reads (FPKM). FPKM is more reflective of the molar amount of a transcript in the original sample than raw counts, as it normalizes the counts for different RNA lengths and for the total number of reads [12]. FPKM is calculated as follows:

$$FPKM = \frac{X_i \times 10^9}{Nl_i} \quad (1)$$

Where X_i is the raw counts for a transcript, l_i is the effective length of the transcript, and N is the total number of counts.

2. Next, we removed non-protein-coding transcripts from the dataset.
3. We removed transcripts that were not meaningfully observed in our dataset by dropping any transcript where $> 70\%$ measurements across all samples were equal to 0.
4. We \log_2 transformed the data
5. We standardized each transcript across all samples, such that the mean for the transcript was equal to zero and the variance of the transcript was equal to one:

$$X'_i = \frac{X_i - \mu_i}{\sigma_i} \quad (2)$$

Attribution Maps using Expected Gradients

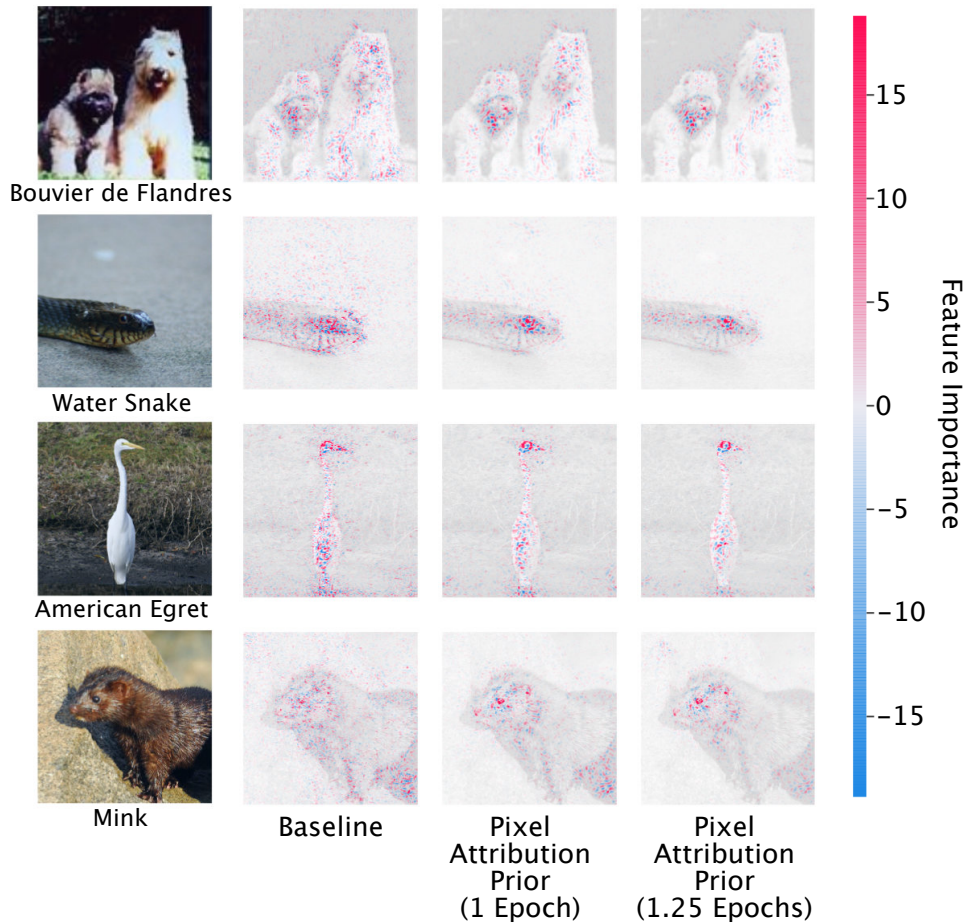


Figure 6: Attribution maps generated by Expected Gradients on the VGG16 architecture before and after fine-tuning using an attribution prior.

where X_i is the expression for a transcript, μ_i is the mean expression of that transcript, and σ_i is the standard deviation of that transcript across all samples.

- Finally, we corrected for batch effects in the measurements using the ComBat tool available in the sva R package [8].

7.2 Train / Validation / Test Set Allocation

To increase the number of samples in our dataset, we opted to use the identity of the drug being tested as a feature, rather than one of a number of possible output tasks in a multi-task prediction. This follows from prior literature on training neural networks to predict drug response [14]. This gave us 30,816 samples (covering 218 patients and 145 anti-cancer drugs). Defining a sample as a drug and a patient, however, meant we had to choose carefully how to stratify samples into our train, validation, and test sets. While it is perfectly legitimate in general to randomly stratify samples into these sets, we wanted to specifically focus on how well our model could learn trends from gene expression data that would generalize to novel patients. Therefore, we stratified samples at a patient-level rather than at the level of individual samples (e.g. no samples from any patient in the test set ever appeared in the training set). We split 20% of the total patients into a test set (6,155 samples), and then split 20% of the training data into a validation set for hyperparameter selection (4,709 samples).

Attribution Maps using Integrated Gradients

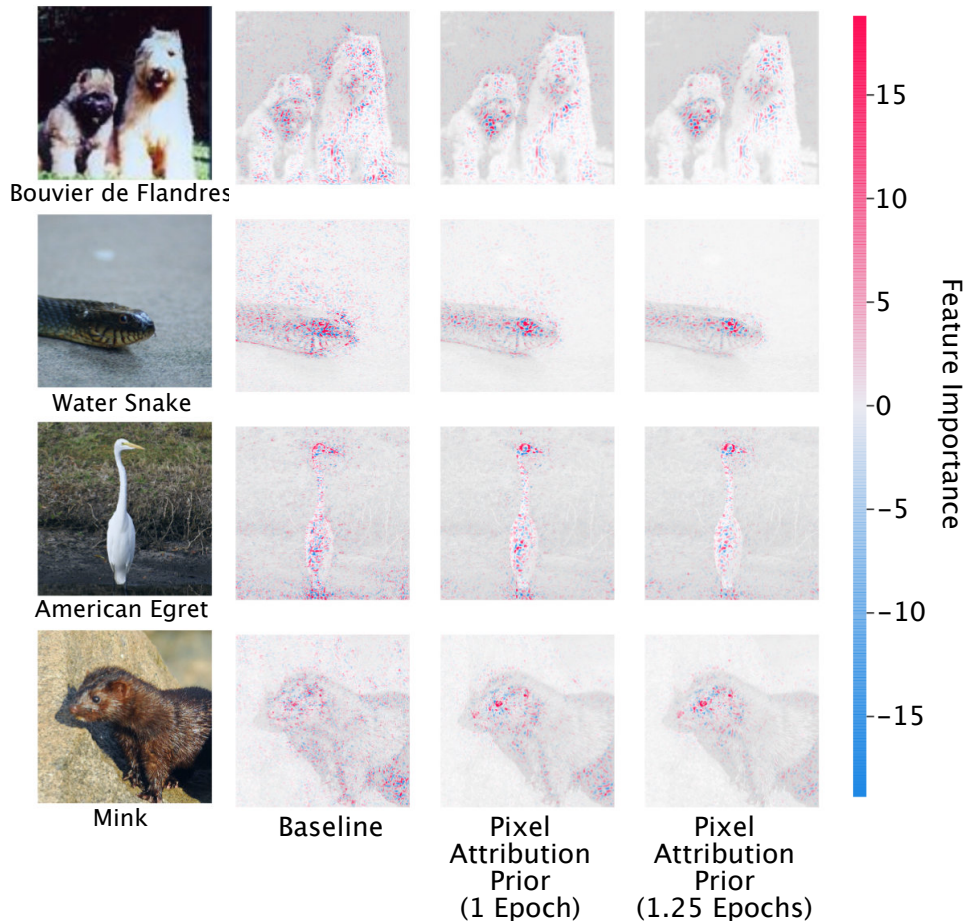


Figure 7: Attribution maps generated by Integrated Gradients on the VGG16 architecture before and after fine-tuning using an attribution prior.

7.3 Model class implementations and hyperparameters tested

LASSO: We used the scikit-learn implementation of the LASSO [24, 13]. We tested a range of α parameters ranging from 10^{-9} to 1, and found that the optimal value for α was 10^{-2} by mean squared error on the validation set.

Graph LASSO: For our Graph LASSO we used the Adam optimizer in TensorFlow [1], with a learning rate of 10^{-5} to optimize the following loss function:

$$\mathcal{L}(w; X, y) = \|Xw - y\|_2^2 + \lambda' \|w\|_1 + \nu' w^T L_G w \quad (3)$$

Where $w \in \mathbb{R}^d$ is the weights vector of our linear model and L_G is the graph laplacian of our HumanBase network [4]. In particular, we downloaded the ‘‘Top Edges’’ version of the hematopoietic stem cell network, which is thresholded to only have non-zero values for pairwise interactions that have a posterior probability greater than 0.1. We used the value of λ' selected as optimal in the regular LASSO model (10^{-2} , corresponds to the α parameter in scikit-learn), and then tuned over a range of ν' values ranging from 10^{-3} to 100. We found that a value of 10 was optimal according to MSE on the validation set.

Neural networks: We tested a variety of hyperparameter settings and network architectures via validation set performance to choose our best neural networks. We tested the following feed-forward network architectures (where each element in a list denotes the size of a hidden layer): [512,256],

Attribution Maps using Gradients x Input

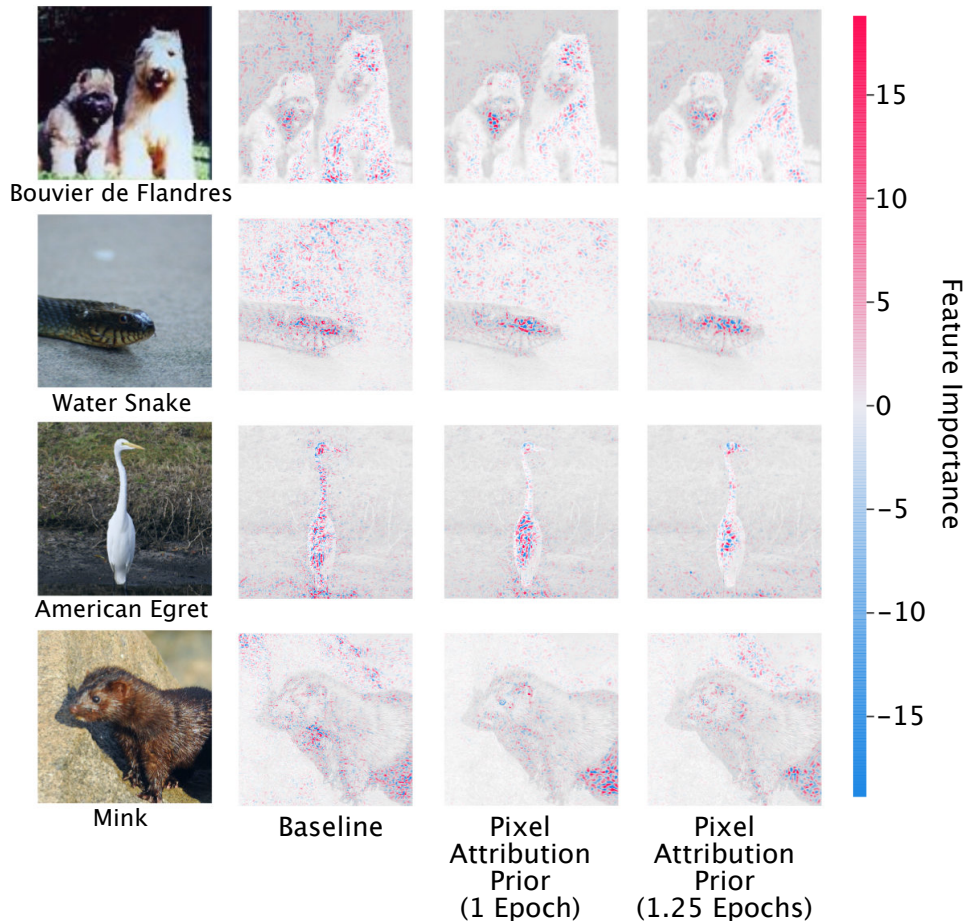


Figure 8: Attribution maps generated by gradient \times input on the VGG16 architecture before and after fine-tuning using an attribution prior.

[256,128], [256,256], and [1000,100]. We tested a range of L1 penalties on all of the weights of the network, from 10^{-7} to 10^{-2} . All models attempted to optimize a least squares loss using the Adam optimizer, with learning rates again selected by hyperparameter tuning from 10^{-5} to 10^{-3} . Finally, we implemented an early stopping parameter of 20 rounds to select the number of epochs of training (training is stopped after no improvement on validation error for 20 epochs, and number of epochs is chosen based on optimal validation set error). We found the optimal architecture (chosen by lowest validation set error) had two hidden layers of size 512 and 256, an L1 penalty on the weights of 10^{-3} and a learning rate of 10^{-5} . We additionally found that 120 was the optimal number of training epochs.

Attribution prior neural networks: To apply our attribution prior to our neural networks, after tuning our networks to the optimal conditions described above, we added extra epochs of fine-tuning where we ran an alternating minimization of the following objectives:

$$\mathcal{L}(\theta; X, y) = \|f_{\theta}(X) - y\|_2^2 + \lambda \|\theta\|_1 \quad (4)$$

$$\mathcal{L}(\theta; X) = \Omega_{graph}(\Phi(\theta, X)) = \nu \bar{\phi}^T L_G \bar{\phi} \quad (5)$$

Following Ross et al. [15], we selected ν to be 100 so that the Ω_{graph} term would be initially equal in magnitude to the least squares and L1 loss term. We found that 5 extra epochs of tuning were optimal by validation set error. We drew $k = 10$ background samples for our attributions.

Attribution Maps using Gradients

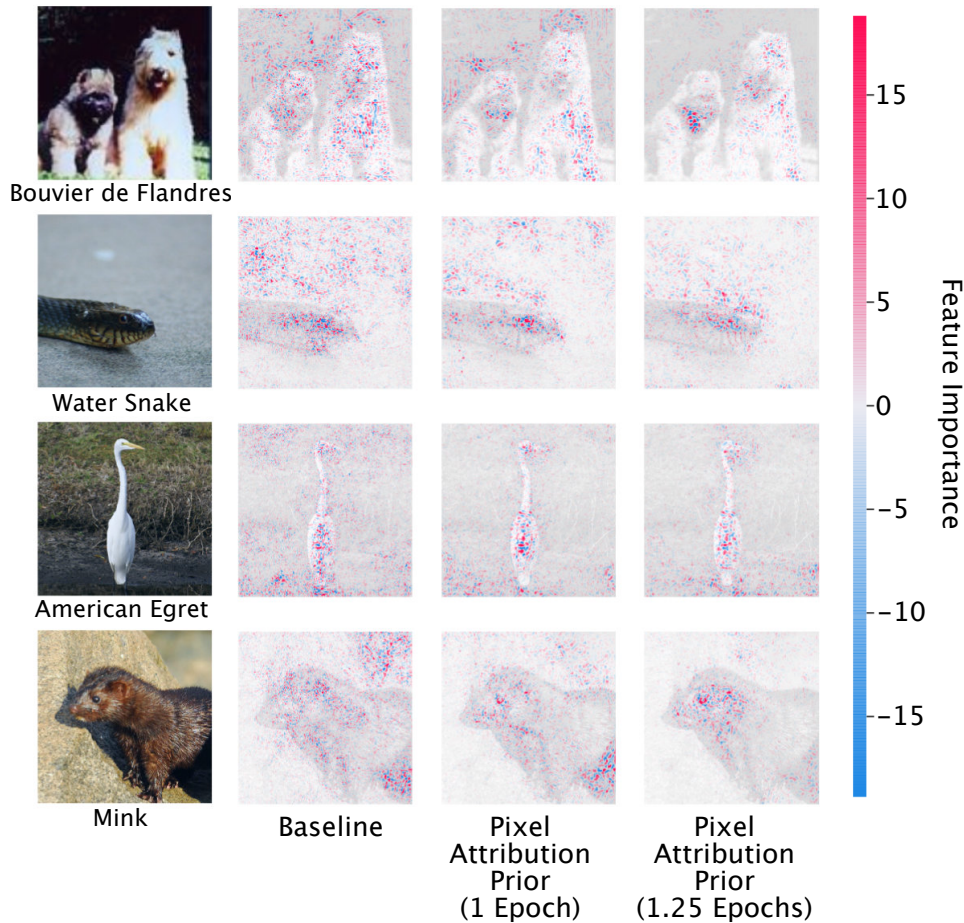


Figure 9: Attribution maps generated by raw input gradients on the VGG16 architecture before and after fine-tuning using an attribution prior.

Graph convolutional networks: We followed the implementation of graph convolution described in Kipf and Welling [7]. The architectures searched were as follows: in every network we first had a single graph convolutional layer (we were limited to one graph convolution layer due to memory constraints on each Nvidia GTX 1080-Ti GPU that we used), followed by two fully connected layers of sizes (512,256), sizes (512,128), or sizes (256,128). We tuned over a wide range of hyperparameters, including L2 penalties on the weights ranging from 10^{-5} to 10^{-2} , L1 penalties on the weights ranging from 10^{-5} to 10^{-2} , learning rates of 10^{-5} to 10^{-3} , and dropout rates ranging from 0.2 to 0.8. We found the optimal hyperparameters based on validation set error were two hidden layers of size 512 and size 256, an L2 penalty on the weights of 10^{-5} , a learning rate of 10^{-5} , and a dropout rate of 0.6. We again used an early stopping parameter and found that 47 epochs was the optimal number.

7.4 Details on experimental results

Looking at the resultant R^2 for prediction, we see that using the graph prior improves predictive performance of a linear model compared to L1-regularization alone (Graph LASSO vs. LASSO). However, we are able to attain a similar degree of predictive performance simply by switching from a linear model to a neural network that does not use the prior graph information at all. Our best performing model was the neural network with graph attribution prior. We use a t -test to compare the R^2 attained from 10 independent retrainings of the neural network to the R^2 attained from

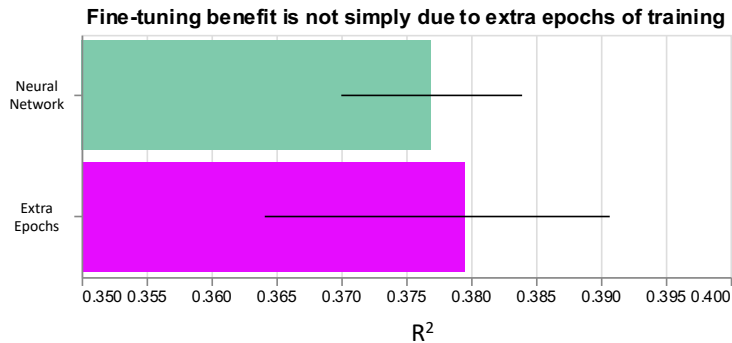


Figure 10: Fine tuning *without* graph prior penalty leads to no significant improvement in model performance.

10 independent retrainings of the attribution prior model and find that predictive performance is significantly higher for the model with the graph attribution prior ($p = 0.004696$).

Since we added graph-regularization to our model by fine-tuning, we wanted to ensure that the improved performance did not simply come from the additional epochs of training *without* the attribution prior. We use a t -test to compare the R^2 attained from 10 independent retrainings of the regular neural network to the R^2 attained from 10 independent retrainings of the neural network with the same number of additional epochs that were optimal when adding the graph penalty (see Figure 10). We found no significant difference between the test error of these models ($p = 0.7565$).

To ensure that the increased performance in the attribution prior model was due to real biological information, we replaced the gene-interaction graph with a randomized graph (symmetric matrix with identical number of non-zero entries to the real graph, but entries placed in random positions). We then compared the R^2 attained from 10 independent retrainings of a neural network with no graph attribution prior to 10 independent retrainings of an neural network regularized with the random graph and found that test error was not significantly different between these two models ($p = 0.5039$). We also compared to graph convolutional neural networks, and found that our network with a graph attribution prior outperformed the graph convolutional neural network ($p = 0.0073$).

To ensure that the models were learning the attribution metric we tried to optimize for, we compared the explanation graph penalty ($\bar{\phi}^T L_G \bar{\phi}$) between the unregularized and regularized models, and found that the graph penalty was on average nearly two orders of magnitude lower in the regularized models (see Figure 12). We also examined the pathways that our top attributed genes were enriched for using Gene Set Enrichment Analysis and found that not only did our graph attribution prior model capture far more significant pathways, it also captured far more AML-relevant pathways (see Figure 11). We defined AML-relevant by a query for the term “AML,” as well as queries for AML-relevant transcription factors.

8 Sparsity experiments

8.1 Data Description

Our data for the sparsity experiments used data from the NHANES I survey [11], and contained 36 variables (expanded to 119 features by one-hot encoding of categorical variables) gathered from 13,000 patients. The measurements include demographic information like age, sex, and BMI, as well as physiological measurements like blood, urine, and vital sign measurements. The prediction task is a binary classification of whether the patient was still alive (1) or not (0) 10 years after data were gathered.

Most important genes for neural network *with attribution prior* come from biologically-relevant pathways

Pathway	FDR q-value
RNA Pol I Promoter Opening	< 10 ⁻²⁸⁰
Amyloids	0.002722
Down-regulated in T Lymphocyte and NK Progenitor cells	0.006435
Down-regulated in normal aging	0.007065
TEL pathway	0.007384
B Cell Lymphoma Cluster 7	0.007601
AML Cluster 9	0.007604
Response to MP470 up	0.007853
Upregulated genes in cells immortalized by HOXA9 and MEIS1	0.008068
AML Cluster 12	0.008163

... +145 more pathways

Most important genes for neural network *without attribution prior* are not significantly enriched for any AML-related pathways

Pathway	FDR q-value
RNA Pol I Promoter Opening	0.001778
Amyloids	0.004001

No additional pathways significant after FDR correction

Figure 11: Top pathways for neural networks with and without attribution priors

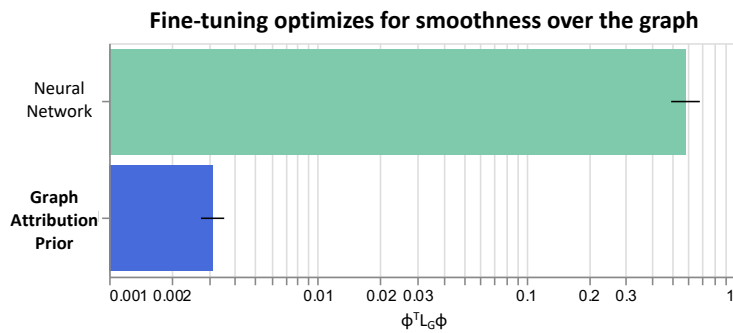


Figure 12: Fine tuning optimizes for the metric we care about: smoothness over the graph

8.2 Data Processing

Data were mean-imputed and standardized so that each feature had 0 mean and unit variance. A fixed train/validation/test split of 7500/2500/3000 patients was used, with all hyperparameter tuning on the validation set and the test set used for final performance plots.

8.3 Model

We trained a range of neural networks to predict survival in the NHANES data. The architecture, nonlinearities, and training rounds were all held constant at values that performed well on an unregularized network, and the type and degree of regularization were varied. All models used ReLU activations and a 2-class softmax output; in addition, all models ran for 1000 batches of 90 samples, or 12 epochs, with an SGD optimizer with learning rate 1.0. 50 references per datum were used for expected gradients attributions in both training and evaluation.

Architecture: We considered a range of architectures including single-hidden-layer 32-node, 128-node, and 512-node networks, two-layer [128,32] and [512,128]-node networks, and a three-layer [512,128,32]-node network; we fixed the [512,128,32] architecture for future experiments.

Regularizers: We tested a large array of regularizers. See 8.4 for details on how optimal regularization strength was found for each regularizer.

- Sparse Attribution Prior - Ω_{sparse} as defined in the main text. After coarse parameter tuning, the optimal regularization strength was found to be $\lambda \in [10^{-2}, 10^0]$, which was the range used for the parameter sweep in the final plots. The best performing model in terms of performance and sparsity, as used in the feature importance distribution plot in the main text, used $\lambda = 1.68 \times 10^{-1}$.
- Mixed L1/Sparse Attribution Prior - Motivated by the observation that the Gini coefficient is normalized and only penalizes the *relative distribution* of global feature importances, we attempted adding an L1 penalty to ensure the attributions also remain small in an absolute sense. This did not result in improvements to overall performance, however. The optimal regularization range for the L1 expected gradients penalty was found to be $\lambda \in [10^1, 10^3]$, and the optimal range for the sparse attribution prior penalty was found to be $\lambda \in [10^{-2}, 10^0]$.
- Sparse Group Lasso - Rather than simply encouraging the weights of the first-layer matrix to be zero, the sparse group lasso also encourages entire columns of the matrix to shrink together by placing an L2 penalty on each column. As in Scardapane et al. [17], we added a weighted sum of column-wise L2 norms to the L1 norm of the entire first-layer matrix, without tuning the relative contribution of the two norms (equal weight on both terms). We optimized the loss directly with SGD. The optimal regularization range was found to be $\lambda \in [10^{-4.7}, 10^{-3.2}] \cup [10^{-1.25}, 10^{-0.75}]$. The best model, as chosen for the main text feature importance distribution plot, used $\lambda = 4.56 \times 10^{-4}$.
- L1 First-Layer - In order to facilitate sparsity, we placed an L1 penalty on the input layer of the network. No regularization was placed on subsequent layers. The optimal regularization range was found to be $\lambda \in [10^{0.35}, 10^{2.35}]$. The best model, as chosen for the main text feature importance distribution plot, used $\lambda = 75.8$.
- L1 All Layers - This penalty places an L1 penalty on all matrix multiplies in the network (not just the first layer). The optimal regularization range was found to be $\lambda \in [10^{0.5}, 10^{2.5}]$.
- L1 Expected Gradients - This penalty penalizes the L1 norm of the vector of global feature attributions, $\bar{\phi}_i$ (analogous to how LASSO penalizes the weight vector in linear regression). The optimal regularization range was found to be $\lambda \in [10^{-0.5}, 10^{0.5}]$.
- L2 First-Layer - This penalty places an L2 penalty on the input layer of the network, with no regularization on subsequent layers. The optimal regularization range was found to be $\lambda \in [10^{-1}, 10^{0.5}]$.
- L2 All Layers - This penalty places an L2 penalty on all matrix multiplies in the network (not just the first layer). The optimal regularization range was found to be $\lambda \in [10^{-2.5}, 10^{-1}]$.

- L2 Expected Gradients - This penalty penalizes the L2 norm of the vector of global feature attributions, $\bar{\phi}_i$ (analogous to how ridge regression penalizes the weight vector in linear models). The optimal regularization range was found to be $\lambda \in [10^{-3}, 10^{-1.5}]$.
- Dropout - This penalty "drops out" a fraction p of nodes at each training batch, but uses all nodes at test time. The optimal regularization range was found to be $\lambda \in [0, 0.1] \cup [0.725, 0.775]$.
- Baseline (Unregularized) - Our baseline model used no regularization.

The Sparse Attribution Prior, Sparse Group Lasso, L1 (All Layers), and Baseline penalties were presented in the main text.

8.4 Hyperparameter tuning:

There was one free parameter to tune for all methods other than the unregularized baseline (no tuning parameter) and the mixed L1/Sparse Attribution Prior model (two parameters - L1 and attribution penalty). We searched all L1 and attribution prior penalties with 130 points sampled on a log scale over $(10^{-10}, 10^3]$. We tuned the dropout probability with 130 points linearly spaced over $(0, 1]$. The mixed L1/Sparse Attribution Prior model was tuned in a 2D grid, with 11 L1 penalties sampled on a log scale over $[10^{-7}, 10^3]$ and 11 attribution prior penalties sampled on a log scale over $[10^{-10}, 10^0]$.

Some penalties, including the sparse attribution prior, mixed, and sparse group lasso penalties, produced NaN outputs for certain regularization settings. We retried several times when NaNs occurred, but if the problem persisted after multiple restarts, the parameter search was truncated and did not search farther out.

We trained 130 models of each type within the "fine" range of penalty strengths determined to be optimal in the "coarse" initial search to illustrate the sparsity/performance tradeoff controlled by regularization strength.

8.5 Results (Full Data)

Performance vs Sparsity Plot: Performance (area under an ROC curve, AUC-ROC) was plotted as a function of sparsity (Gini coefficient) for all models. Figure 14 shows sparsity and validation performance for the coarse initial parameter sweep, as well as sparsity and test performance for the fine sweep for all models. The third image in the figure is a zoomed version to provide more detail on the best-performing models. The three model classes shown in the main text were those that were capable of both high accuracy and high sparsity; none of the other models outperformed L1, SGL, or sparse attribution prior models in terms of ROC-AUC, and most were substantially less sparse.

Feature Importance Distribution Plot: The distribution of feature importances was plotted in the main text as a Lorenz curve: for each model, the features were sorted by global attribution value $\bar{\phi}_i$, and the cumulative normalized value of the lowest q features was plotted, from 0 at $q = 0$ to 1 at $q = p$. A lower area under the curve indicates more features have relatively small attribution values, indicating the model is sparser.

Feature Importance Summary: We also show summaries of the attributions for the top 20 features in each model in Figure 15. Each row shows the distribution of attributions for a given feature in a given model across the entire training set. The color indicates the value of the feature; for example, red dots to the left in the "age" row indicates that high age reduces probability of survival at ten years. The sparse attribution prior places substantially less weight than other models on features that are not in the top 6 by importance.

8.6 Small-Data Experiments

The small-data experiments used small subsamples of the training and validation sets to test ability of sparse models to learn with limited data. The data description and data processing were identical to the larger-data experiments, except that the train/validation/test split was 100/100/2500. 10 datasets were created using random subsets of the training and validation sets. Models and architectures were identical to the larger-data experiments, but the batch size was the entire dataset (100 points), and only 20 epochs/batches were run (increasing to 100 epochs/batches did not improve performance of

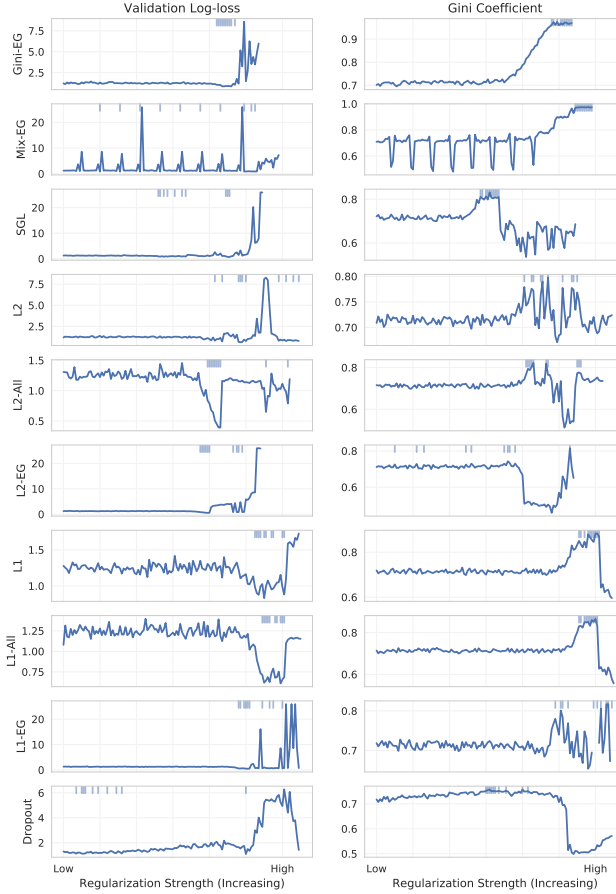


Figure 13: Validation loss and gini coefficient as a function of regularization strength for all models. Blue hashes at top of each plot indicate the 10 points with lowest loss, and highest gini coefficient, in the respective columns. The areas with lowest log-loss were used to determine the range of fine parameters to sweep for the final result.

unregularized models). The best-performing models from the larger-data experiment were tested, including Sparse Attribution Prior, Sparse Group Lasso, L1 (All Layers), and Baseline penalties. Hyperparameters were tuned over the same log scales as in the larger-data experiment, but were not fine-tuned across a smaller range. Because the small-data models could be trained faster, we continued the parameter search if repeated NaNs occurred, skipping a parameter setting if NaNs occurred on 10 random restarts. The best hyperparameters for each model on a given dataset were selected using validation data, and 100 new models were trained with these parameters. The top 10 resulting models (by validation error) were evaluated on the full test data to yield final results. Differences in performance were evaluated using a paired-samples T -test.

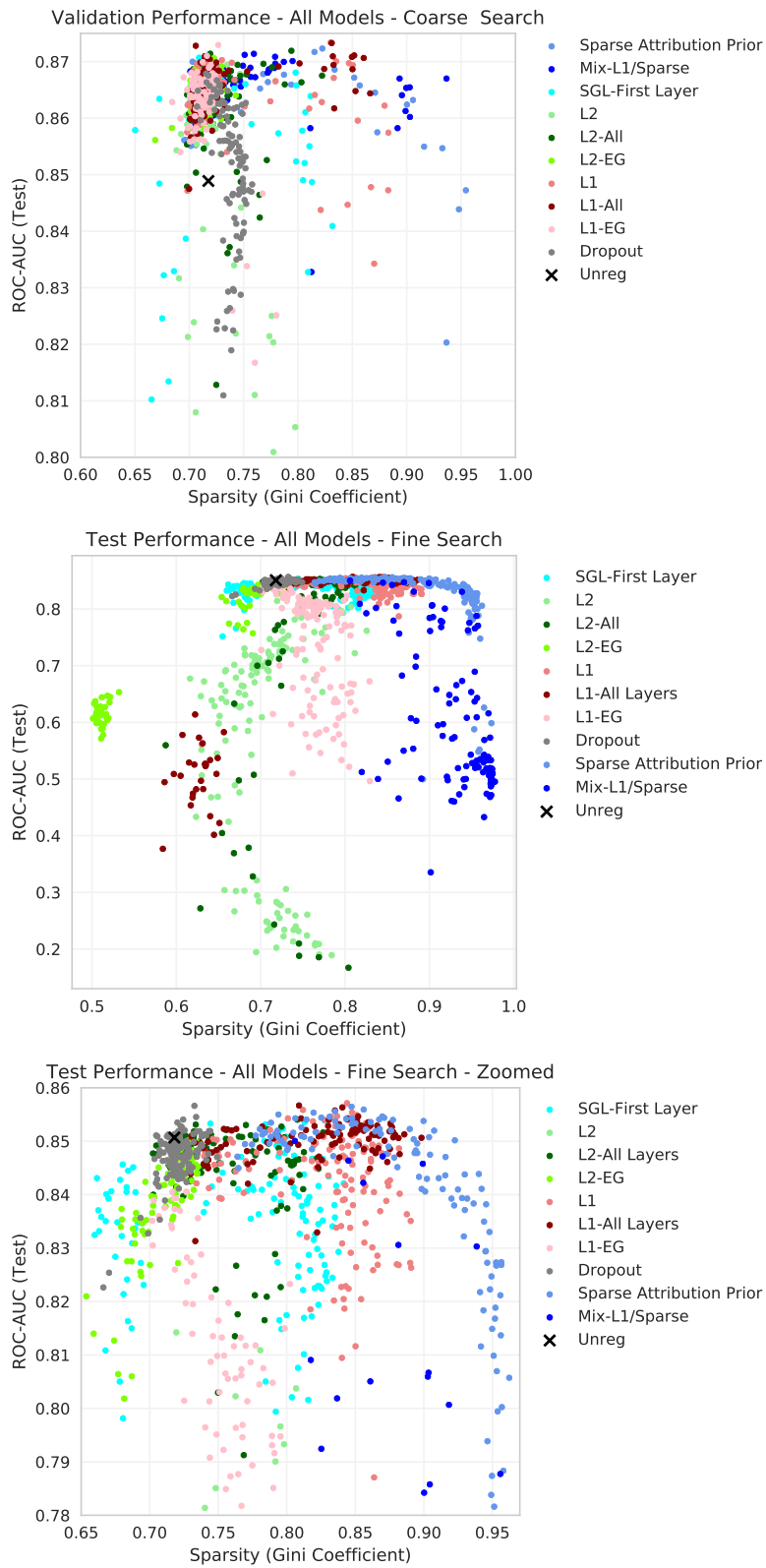


Figure 14: Sparsity vs performance plot for all models.

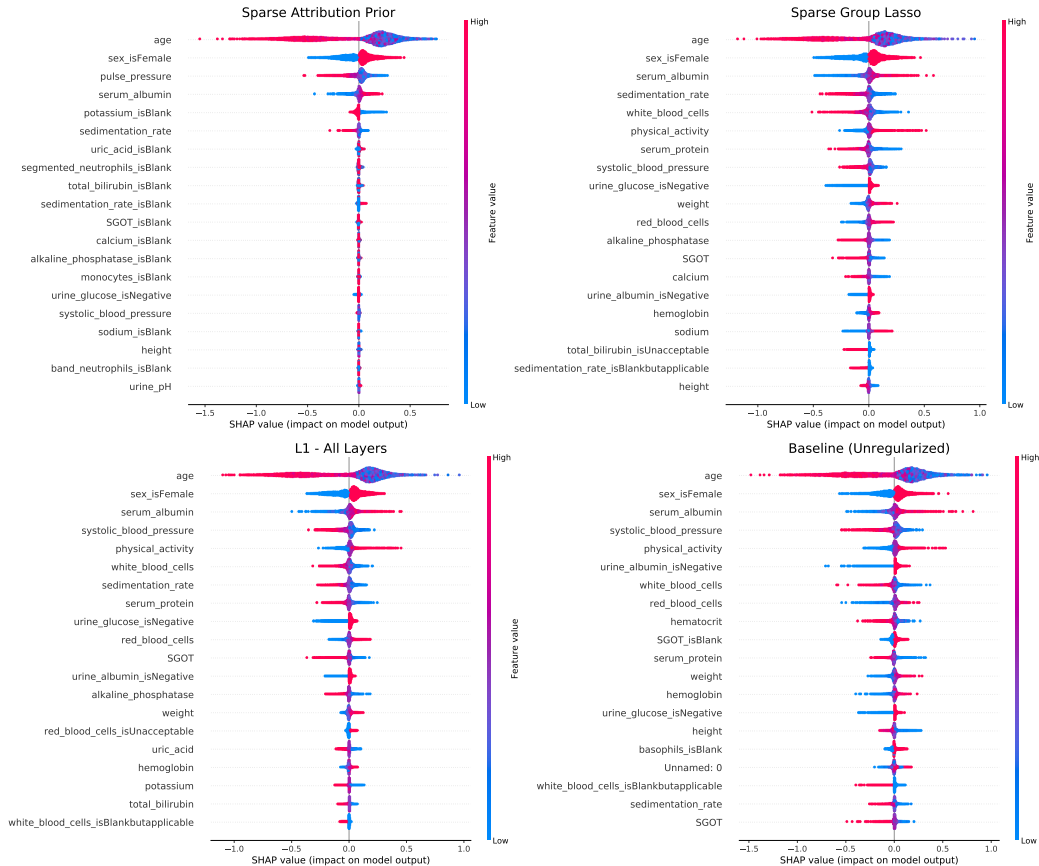


Figure 15: Summary of feature attributions for top 20 features from each model (best model from each class chosen as described in the main text). Note that SGL and L1 are somewhat sparser than unregularized models, but the sparse attribution prior assigns dramatically less importance to all but the top few features.

References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.
- [2] J. M. Bardsley. Laplace-distributed increments, the laplace prior, and edge-preserving regularization. *J. Inverse Ill-Posed Probl*, 2012.
- [3] H. Drucker and Y. Le Cun. Improving generalization performance using double backpropagation. *IEEE Transactions on Neural Networks*, 3(6):991–997, 1992.
- [4] C. S. Greene, A. Krishnan, A. K. Wong, E. Ricciotti, R. A. Zelaya, D. S. Himmelstein, R. Zhang, B. M. Hartmann, E. Zaslavsky, S. C. Sealfon, et al. Understanding multicellular function and disease with human tissue-specific networks. *Nature genetics*, 47(6):569, 2015.
- [5] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019.
- [6] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [7] T. N. Kipf and M. Welling. Semi-Supervised Classification with Graph Convolutional Networks. *CoRR*, abs/1609.0, 2016. URL <http://arxiv.org/abs/1609.02907>.

- [8] J. T. Leek and J. D. Storey. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLOS Genetics*, 3(9):1–12, 09 2007. doi: 10.1371/journal.pgen.0030161. URL <https://doi.org/10.1371/journal.pgen.0030161>.
- [9] Y. Lou, T. Zeng, S. Osher, and J. Xin. A weighted difference of anisotropic and isotropic total variation model for image processing. *SIAM Journal on Imaging Sciences*, 8(3):1798–1823, 2015.
- [10] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee. Explainable ai for trees: From local explanations to global understanding, 2019.
- [11] H. W. Miller. Plan and operation of the health and nutrition examination survey, united states, 1971-1973. DHEW publication no.(PHS)-Dept. of Health, Education, and Welfare (USA), 1973.
- [12] A. Mortazavi, B. A. Williams, K. McCue, L. Schaeffer, and B. Wold. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nature Methods*, 5:621, may 2008. URL <https://doi.org/10.1038/nmeth.1226><http://10.0.4.14/nmeth.1226><https://www.nature.com/articles/nmeth.1226#supplementary-information>.
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [14] K. Preuer, R. P. I. Lewis, S. Hochreiter, A. Bender, K. C. Bulusu, and G. Klambauer. Deep-Synergy: predicting anti-cancer drug synergy with Deep Learning. *Bioinformatics*, 34(9):1538–1546, 2018. doi: 10.1093/bioinformatics/btx806. URL <http://dx.doi.org/10.1093/bioinformatics/btx806>.
- [15] A. S. Ross, M. C. Hughes, and F. Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. *arXiv preprint arXiv:1703.03717*, 2017.
- [16] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [17] S. Scardapane, D. Comminiello, A. Hussain, and A. Uncini. Group sparse regularization for deep neural networks. *Neurocomputing*, 241:81–89, 2017.
- [18] Y. Shi and Q. Chang. Efficient algorithm for isotropic and anisotropic total variation deblurring and denoising. *Journal of Applied Mathematics*, 2013, 2013.
- [19] N. Silberman and S. Guadarrama. Tensorflow-slim image classification model library. 2016.
- [20] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [21] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [22] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. *JMLR. org*, 2017.
- [23] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [24] R. Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. ISSN 00359246. URL <http://www.jstor.org/stable/2346178>.

- [25] F. Yu, Z. Xu, Y. Wang, C. Liu, and X. Chen. Towards robust training of neural networks by regularizing adversarial gradients. [arXiv preprint arXiv:1805.09370](#), 2018.
- [26] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. [arXiv preprint arXiv:1710.09412](#), 2017.