
xGEMs: Generating Exemplars to Explain Black-Box Models

Shalmali Joshi
UT Austin
shalmali@utexas.edu

Oluwasanmi Koyejo
UIUC
sanmi@illinois.edu

Been Kim
Google Brain
beenkim@google.com

Joydeep Ghosh
UT Austin
jghosh@utexas.edu

Abstract

This work proposes **xGEMs**: or *manifold guided exemplars*, a framework to understand black-box classifier behavior by exploring the landscape of the underlying data manifold as data points cross decision boundaries. To do so, we train an unsupervised implicit generative model – treated as a proxy to the data manifold. We summarize black-box model behavior quantitatively by perturbing data samples along the manifold. We demonstrate **xGEMs**' ability to detect and quantify bias in model learning and also for understanding the changes in model behavior as training progresses.

1 Introduction

Machine learning algorithms have become widely deployed in domains beyond web based recommendation systems, like the criminal justice system [3], clinical healthcare [6] etc. For instance, risk assessment tools like COMPAS [3] produce learned recidivism scores to consequently determine the amount of pre-trial bail and detention. Similarly, medical interventions can impact health outcomes for patients, making institutions liable to provide explanations for their decisions. This has motivated regulatory agencies like the EU Parliament¹ to codify a right to data protection and “obtain an explanation of the decision reached using such automated systems²”.

Systems that provide satisfactory explanations for the decisions of such learning algorithms have until recently been few and far between. It is challenging to characterize the specific nature of explainability mechanisms given their complexity and lack of consensus on the nature and sufficiency of such explanations [11, 27]. The problem is often compounded due to multiple levels of abstraction required to provide such explanations [29]. For instance, system level explanations as required by regulatory bodies are different from an abstraction that would assist practitioners of machine learning. This work focuses on providing explanations for low level understanding of model behavior, albeit at an abstraction beyond performance metrics. Such a suite of explanations not only help improve

¹in collaboration with the European Commission and the Council of the European Union

²<https://www.privacy-regulation.eu/en/r71.htm>

understanding of opaque models³ [19, 21] but can also uncover biases (inherent in the data) that models pick up on e.g. learned gender and racial biases [5].

In this work, we posit that there need not be an inherent tradeoff between model performance and explainability, as is generally assumed (and found in [22, 16, 20]). We propose an explainability tool that probes a supervised black-box model along the data manifold for explanations via examples and/or summaries. Demonstrating model behavior via examples is known to be beneficial for improving and understanding the decision making process [1]. Navigating the data manifold allows us to explore black-box behavior in different regions of the manifold. The proposed method can be utilized as a diagnostic tool to analyze training progression, compare classifier performance, and/or uncover inherent biases the classifier may have learned.

2 Related Work

Most closely related works to our approach are those that provide explanations by sub-selecting meaningful samples and/or semantically relevant features (like super-pixels) that highlight undesirable model behavior [13, 22]. Most of these methods require the selected samples to be part of training/test dataset. This means that if the training/test set did not include the instance that best explains a specific decision, we would have to settle for a suboptimal choice. Our method aims to relax this constraint by generating new examples that are better suited for this purpose. In terms of generating examples, adversarial criticisms [35] and the class of generative networks like GANs are relevant approaches. Specifically, [35] use the adversarial attack paradigm as a means to select examples from existing training data to explain model behavior, similar to [22]. However note that the goal of generating adversarial examples and our explanations are fundamentally different. The primary goal of adversarial examples is to focus on exploiting the worst case confounding scenario given a decision boundary, while our work focuses on generating an example that lies on the data manifold as it crosses a decision boundary. See Figure 1a for a more intuitive explanation. We posit that it is important to uncover classifier behavior when data points are constrained to the data manifold. Such data instances are more 'realistic' and likely to be created by the underlying phenomenon that led to the training data. They provide an alternative method to probe a black-box, specially in non-adversarial settings. They also characterize the residual vulnerabilities of a model that defends itself against adversarial attacks by detecting directed "noise" that is orthogonal to the manifold of the data or of an associated latent space.

We position our work as a diagnostic framework for understanding model behavior at an abstraction that may be most useful to a data science practitioner and/or a machine learning expert. However, as suggested before, explainable models focus on different notions of explainability. For example, Koh and Liang [25] use influence functions, motivated by robust statistics Cook and Weisberg [7] to determine importance of each training sample for model predictions. Li et al. [26], Selvaraju et al. [31] focus on understanding the workings of different layers of a deep network and studying saliency maps for feature attribution [33, 34, 36]. Saliency methods, while powerful, can be demonstrated to be unreliable without stronger conditions over the saliency model [23, 2]. Other paradigms of explainable models focus on locally approximating complex models using a simpler functional form to approximate the (local) decision boundary. For instance, LIME based approaches [30, 32, 4] locally approximate complex models with linear fits. Decision Trees are also considered more explainable if they are not too large. These approaches inherently assume a tradeoff between model performance and explainability, as less complex model classes tend to be empirically sub-par in performance relative to the success of the target black-box models they endeavor to explain. The **xGEMs** framework, however, does not rely on local approximations to provide explanations or assume such a trade-off.

We summarize our key contributions as follows: 1. We introduce **xGEMs**, a framework for explaining supervised black-box models via examples generated along the underlying data manifold. 2. We demonstrate the utility of **xGEMs** in (a) detecting bias in learned models, (b) characterizing the probabilistic decision manifold w.r.t. examples, and (c) facilitating model comparison beyond standard performance metrics.

³<https://distill.pub/2018/building-blocks/>

Algorithm 1 Find (\mathbf{x}^*, y^*) -xGEM

Input: $(\mathbf{x}^*, y^*) \in \mathbb{R}^d \times \{-1, 1\}, y_{tar}, \mathcal{G}_\theta, \mathcal{F}_\psi, f_\phi, \lambda, \eta > 0$

```

Initialize  $\mathbf{z} = \mathcal{F}_\psi(\mathbf{x}^*)$ 
while Not converged do
     $\tilde{\mathbf{z}} \leftarrow \tilde{\mathbf{z}} + \eta \nabla_{\tilde{\mathbf{z}}} (\mathcal{L}(\mathbf{x}^*, \mathcal{G}_\theta(\tilde{\mathbf{z}})) + \lambda \ell(f_\phi(\mathcal{G}(\tilde{\mathbf{z}})), y_{tar}))$ 
     $\tilde{\mathbf{x}} = \mathcal{G}_\theta(\tilde{\mathbf{z}})$ 
Return  $\tilde{\mathbf{x}}$ 

```

3 Background

Implicit Generative Models can be described as stochastic procedures that generate samples (denoted by the random variable $\mathbf{x} \in \mathbb{X}^d$) from the data distribution $p(\mathbf{x})$ without explicitly parameterizing $p(\mathbf{x})$. The two most significant types are the Variational Auto-Encoders (VAEs) [24] and Generative Adversarial Networks (GANs) [14]. Implicit generative models generally assume an underlying latent dimension $\mathbf{z} \in \mathbb{R}^k$ that is mapped to the ambient data domain $\mathbf{x} \in \mathbb{R}^d$ using a deterministic function \mathcal{G}_θ parametrized by θ , usually as a deep neural network. The primary difference between GANs and VAEs is the training mechanism employed to learn function \mathcal{G}_θ . GANs employ an adversarial framework by employing a discriminator that tries to classify generated samples from the deterministic function versus original samples and VAEs maximize an approximation to the data likelihood. The approximation thus obtained has an encoder-decoder structure of conventional autoencoders [9]. One can obtain a latent representation of any data sample within the latent embedding using the trained encoder network. While GANs do not train an associated encoder, recent advances in adversarially learned inference like BiGANs [12, 10] can be utilized to obtain the latent embedding. In this work, we assume access to an implicit generative model that allows us to obtain the latent embedding of a data point.

Let $\mathcal{F}_\psi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ (parametrized by ψ) be the inverse mapping function that provides the latent representation for a given data sample. Let $\mathcal{L} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ be the analogous loss function such that for a given data sample $\tilde{\mathbf{x}}$:

$$\tilde{\mathbf{z}} = \arg \min_{\mathbf{z}} \mathcal{L}(\tilde{\mathbf{x}}, \mathcal{G}_\theta(\mathbf{z})) \triangleq \mathcal{F}_\psi(\tilde{\mathbf{x}}) \quad (1)$$

Examples of \mathcal{F}_ψ are the encoder in a VAE, or an inference network in a BiGAN. An appropriate distance function in the data domain can be used as the loss \mathcal{L} .

Without loss of generality, we assume that we would like to provide explanations for a binary classifier. Let $y \in \{-1, 1\}$ be the target label. Let $f_\phi : \mathbb{R}^d \rightarrow \{-1, 1\}$ be the target black-box classifier to be ‘explained’ and $\ell(f_\phi(\mathbf{x}), y)$ be the loss function used to train the black-box classifier.

Adversarial Criticisms Adversarial criticisms to explain black-box classifiers look for perturbations $\delta_{\mathbf{x}}$ to data samples \mathbf{x} such that the perturbations maximize the loss $\ell(f_\phi(\mathbf{x} + \delta_{\mathbf{x}}), y)$ or change the predicted label. These perturbations are invisible to the human eye. That is, if $\tilde{\mathbf{x}}$ is the target adversarial sample, an adversarial attack solves a Taylor approximation to the following:

$$\tilde{\mathbf{x}} = \arg \max_{\tilde{\mathbf{x}}: \|\tilde{\mathbf{x}} - \mathbf{x}\|_p < \epsilon} \ell(f_\phi(\tilde{\mathbf{x}}), y) \quad (2)$$

We now characterize the proposed model and detail the kinds of explanations it can provide.

4 Generating xGEMs

To provide explanations via examples over more *naturalistic* perturbations, we introduce a new set of examples, called *manifold guided examples* or **xGEMs**. First, we train an implicit generative model \mathcal{G}_θ and an encoder network \mathcal{F}_ψ .

$$\tilde{\mathbf{x}} = \mathcal{G}_\theta(\arg \min_{\mathbf{z} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}^*, \mathcal{G}_\theta(\mathbf{z})) + \lambda \ell(f_\phi(\mathcal{G}(\mathbf{z})), y_{tar})) \quad (3)$$

A manifold guided example is defined w.r.t. a given data sample \mathbf{x}^* .

Definition 1 (\mathbf{x}^*, y^* -xGEM). An **xGEM** corresponding to a data point (\mathbf{x}^*, y^*) and a target label $y_{tar} \neq y^*$, refers to the solution of Equation (3) for a fixed and known $\lambda > 0$. The **xGEM** is denoted by $\tilde{\mathbf{x}}$.

We propose Algorithm 1 to estimate a manifold guided example **xGEM** for any data point \mathbf{x}^* . Intuitively, for a point \mathbf{x}^* , we determine its latent representation using \mathcal{F}_ψ . This allows us to explain model behavior from a common latent representation across all black-boxes. To find realistic perturbations to this point, along the data manifold, we traverse the latent space of the generator \mathcal{G}_θ (our proxy for the data manifold) until the label switches to the desired target label y_{tar} . The desired *manifold guided example* or **xGEM** is the sample generated at the switch point in the latent embedding.

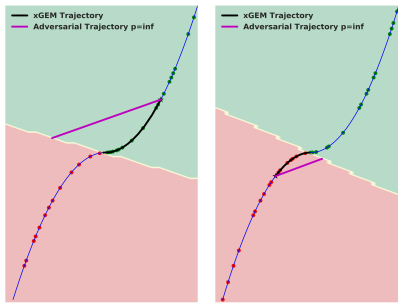
We empirically highlight the benefits of the discovering manifold guided examples in different contexts and abstractions that provide insights into model behavior.

5 Explanations using xGEMs

We first use a simple setting with simulated data to highlight the differences between the proposed explanation tool compared to criticisms and prototypes derived from adversarial attacks [35].

5.1 An alternative view to Adversarial Criticisms

Figure 1a demonstrates a linear decision boundary trained on data with ambient dimension equal to 2. The one-dimensional data manifold is parabolic as shown by the blue curve. The green points are in class 1 and red points are samples belonging to class 2. The figure illustrates manifold guided examples as well as the trajectory taken by the gradient steps of Algorithm 1. The trajectory to generate an adversarial criticism stems from Equation (2). A generative model maps from a 1d latent dimension to the data manifold shown by the blue curve. A single layer (softmax) neural network with output dimension=2 is trained on points sampled from this manifold (the yellow decision boundary separates the two classes – regions marked by the pink and green regions). As demonstrated by the figure, navigating along the latent dimension of the generator encourages the **xGEM** trajectory to be constrained along the data manifold, while adversarial criticisms may lie well outside the manifold. Thus *manifold guided examples* offer alternative view of classifier behavior via examples.



(a)

Figure 1

xGEMs versus *Adversarial* criticisms [35], for a parabolic manifold (shown in blue). Green points belong to class 1 and red points to class 2. The black trajectories in all figures are gradient steps taken by Algorithm 1 while the magenta trajectories correspond to adversarial trajectories determined by Equation 2 with $p = \infty$. Note that all decision boundaries in Figures (a) and (b) separate the data. The decision boundary is trained by optimizing a softmax regression using the cross-entropy loss function.

We defer examples of **xGEMs** evaluated for the MNIST dataset to the Appendix in the interest of space.

5.2 Towards automated bias detection

We now demonstrate the utility of generating manifold guided examples to detect if a target classifier is confounded w.r.t. a given attribute of interest. In particular, we wish to determine whether a black-box is differentiating among the target labels using spurious correlations in the data. For instance, a classifier trained to determine the best medical intervention may be relying on attributes like gender to determine best treatment. It is desirable to have an automated mechanism to detect such behavior. We say that a classifier is confounded with an attribute of interest a if the attribute a substantially influences the black-box’s predictions. We make this concrete in the context of our framework below.

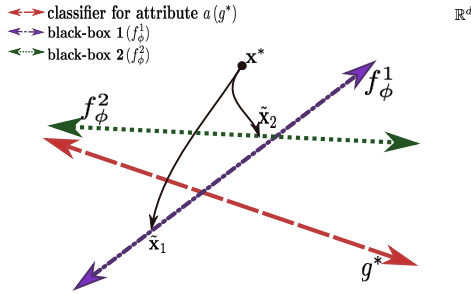


Figure 2: Example of bias detection. Target black-boxes: f_ϕ^1 and f_ϕ^2 . g^* classifies points w.r.t. a . \tilde{x}_1 and \tilde{x}_2 are **xGEMs** corresponding to x^* for f_ϕ^1 and f_ϕ^2 resp. \tilde{x}_2 's attribute prediction (w.r.t g^*) is the same as that of x^* while that of \tilde{x}_1 is different. Thus we say that f_ϕ^1 is biased w.r.t. attribute a for sample x^* .

Without loss of generality let $a \in \{-1, 1\}$ be the (potentially protected) binary attribute of interest. We wish to examine whether the target classifier f_ϕ is biased/confounded by a . Intuitively, we hope that attribute a of an **xGEM** should be the same as that of the original point. In order to detect this, we assume there exists an oracle $g^* : \mathbb{R}^d \rightarrow \{-1, 1\}$ that perfectly classifies the confounding attribute a when considered as the dependent variable, based on the other (d) independent variables. Additionally, we assume that g^* is not confounded by the target label of the black-box y and is not used by g^* to predict a . Let $\mathbb{R}^d \times \{-1, 1\} \times \{-1, 1\} \supset \mathcal{D} \triangleq \{(\mathbf{x}_i, y_i, a_i), i \in [N]\}$ be the training data where i indexes a given point. Let \tilde{x}_i be the **xGEM** of x_i w.r.t. f_ϕ as returned by Algorithm 1. We argue that classifier f_ϕ is confounded by the attribute a if equation (4) holds for a given $\delta > 0$.

$$\frac{E_{\mathcal{D}}[\mathbb{1}(g^*(\tilde{x}) \neq a)]}{|\mathcal{D}|} > \delta \quad (4)$$

In practice, access to a perfect oracle g^* is infeasible or prohibitively expensive. In some cases, such a classifier may be provided by regulatory bodies, thereby adhering to predetermined criterion as to what accounts for a *reliable* proxy oracle. For this case study, we assume it is sufficient that the proxy oracle has the same false positive and false negative error rates w.r.t. the target label, which is a fairness condition known as the Equalized Odds Criterion [17]. To demonstrate our algorithm, we assume access to a proxy oracle $\hat{g} : \mathbb{R}^d \rightarrow \{-1, 1\}$ that satisfies the following conditions, given a $0.5 \ll \tau < 1$:

- (i) $E_{\mathcal{D}}[\mathbb{1}(\hat{g} == a)] > \tau$ (ii) \hat{g} satisfies the Equalized Odds [17] criterion w.r.t. the target label y .

Black-box Classifier	Accuracy	Confounding metric
f_ϕ^1	0.9933	0.1704
f_ϕ^2	0.9155	0.4323

Table 2: Confounding metric

Black-box	Target label	
	Black Hair	Blond Hair
f_ϕ^1	Male:0.4550	Male:0.1432
	Female:0.0159	Female:0.0484
	Overall:0.2430	Overall:0.0539
f_ϕ^2	Male:0.7716	Male:0.1475
	Female:0.0045	Female:0.5024
	Overall:0.4012	Overall:0.4821

Table 3: Confounding metric by gender

Note that while we consider \hat{g} as an inexpensive proxy for g^* , we prescribe that the experiment be carried out with g^* . Figure 2 demonstrates how such confounding could be detected, as well as used for model comparison w.r.t. their biases. As shown in the figure, f_ϕ^1 and f_ϕ^2 are the classification boundaries of two black-box models classifying a target label of interest. g^* is a classifier that

Attribute Classifier	(a)	Target black-box label	
		Black Hair	Blond Hair
\hat{g} (orig)		FP:0.003	FP:0.000
		FN:0.002	FN:0.018
		Acc: 0.997	Acc:0.999
\hat{g} (recalibrated)		FP:0.003	FP:0.003
		FN:0.018	FN:0.018
		Acc:0.989	Acc:0.996

Table 1: Recalibrated Gender Classifier.

classifies the data according to attribute a . Consider the sample \mathbf{x}^* and let $\tilde{\mathbf{x}}_1$ and $\tilde{\mathbf{x}}_2$ be the manifold guided examples of \mathbf{x}^* corresponding to classifiers f_ϕ^1 and f_ϕ^2 respectively. As shown in the figure, the attribute a of the **xGEM** $\tilde{\mathbf{x}}_1$ is different from that of \mathbf{x}^* while that of $\tilde{\mathbf{x}}_2$ is not. We conclude that a black-box f_ϕ^1 is confounded if the fraction of points whose manifold guided examples or **xGEMs** that change attribute a is greater than δ . Thus an empirical estimate of Equation (4) gives a metric that can quantify the amount of confounding in a given black-box, while also allowing to compare different black-boxes w.r.t. the target attribute a .

We evaluate our framework for confounding detection in facial images using the CelebA [28] dataset. The target black-box classifier predicts the binary facial attribute – hair color (black or blond). We determine whether or not the black-box is confounded with gender. We restrict to two genders, male and female, based on annotations available in CelebA. In particular, \hat{g} is a ResNet model [18]⁴ that classifies celebA faces by gender. \hat{g} is recalibrated to satisfy the two conditions mentioned earlier. Details of \hat{g} 's performance and recalibration are provided in Table 1.

Two ResNet models f_ϕ^1 and f_ϕ^2 are trained to detect the hair color attribute (black hair vs blond hair) using two different datasets. f_ϕ^1 is trained on all face samples with either black or blond hair whereas f_ϕ^2 is trained such that all black hair samples are male while blond haired samples are all female. Table 2 gives the overall validation accuracy of both classifiers. Note that the validation set used for f_ϕ^1 and f_ϕ^2 are the same.

Table 2 also shows the fraction of samples whose manifold guided examples' predicted attribute a (in this case gender) is different from the original training sample w.r.t. \hat{g} . The fraction of confounded samples is clearly much larger for the classifier trained on a biased dataset as determined by the proxy oracle \hat{g} . Additionally, Table 3 suggests a 10-fold increase in the fraction of confounding for blond haired females with the biased classifier f_ϕ^2 . Notice the decrease in the amount of confounding for black haired females while a general increase in confounding for all black haired faces. As an aside, the biased model f_ϕ^2 also changes the background more than hair color in order to change the hair color label (see Figure 3). This suggests that quantifying such confounding using manifold guided examples allows us to characterize biases w.r.t. any attribute of interest.

Figure 3 shows a few examples of such confounded images for the two black-boxes. In particular, we show examples where the black-box trained on biased data for hair color classification changes gender of the sample as it crosses the decision boundary whereas the black-box trained on unbiased data does not⁵.

5.3 Case Study: Model Assessment beyond performance metrics

An important aspect of black-box analyses is to study the progression of training complex models. Specifically, observing manifold guided examples allows us to consider model behavior in the following aspects: 1) Discerning shifts in features relied on by the black-box to differentiate between classes during training. 2) Characterizing the probabilistic manifolds of manifold guided examples as training progresses and its relation to calibration of complex networks [8]. 3) Qualitative trade-offs and/or mistakes made by the classifier for prototypical examples.

Reliability Diagrams have been used as a summary statistic to evaluate model calibration [8] that aims to study whether the confidence of a prediction matches the ground truth likelihood of the prediction. It has been observed that while model performance has improved substantially in recent years because of deep networks, such models are typically more prone to mis-calibration [15]. We provide a complementary statistic to Reliability Diagrams to assist model assessment/comparison.

For this study we train two deep networks \hat{f}_ϕ^1 (a ResNet model) and \hat{f}_χ^2 (a four layer CNN with local response normalization (lrn)⁶) with CelebA face images for the hair color (black/blond) binary classification task. For a given face, we evaluate the corresponding **xGEM** at multiple incremental training steps. We plot the confidence of labeling a point to have black hair with respect to the distance

⁴https://github.com/tensorflow/models/tree/master/tutorials/image/cifar10_estimator

⁵All qualitative figures were chosen based on the confidence of the prediction from the black-box and confidence of the reconstructed image

⁶<https://github.com/tensorflow/models/tree/master/tutorials/image/cifar10>



Figure 3: We test whether ResNet models f_ϕ^1 and f_ϕ^2 , both trained to detect hair color but on different data distributions are confounded with gender. Two samples for classifiers f_ϕ^1 (first sub row) and f_ϕ^2 (second sub row) are shown. The leftmost image is the original figure, followed by its reconstruction from the encoder F_ψ . Reconstructions are plotted as Algorithm 1 (with $\lambda = 0.01$) progresses toward crossing the decision boundary. The red bar indicates change in hair color label indicated at the top of each image along with the confidence of prediction. The label at the bottom indicates gender as predicted by \hat{y} . For both samples, classifier f_ϕ^1 , trained on biased data changes the gender (1st and 3rd rows) while crossing the decision boundary whereas the other black-box does not.

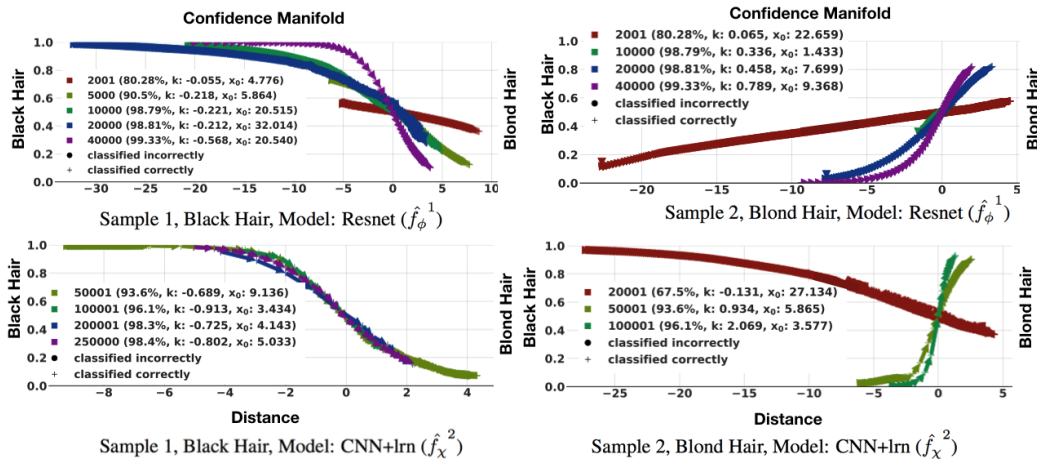


Figure 4: Confidence manifolds for a few data samples for black-box models 1 and 2.

of the original reconstruction and its **xGEM** including all intermediate points from the decision boundary (called ‘confidence manifold’). Thus, all samples originally labeled black should have high confidence of being labeled and the confidence decreases as the sample crosses the decision boundary (vice-versa for blond haired faces). Figure 4 shows the confidence manifolds for two samples (one in each column).

The top and bottom rows represent the manifolds obtained during training for model 1 (\hat{f}_ϕ^1) and model 2 (\hat{f}_χ^2) respectively. Sample 1(column 1) is a face with black hair while Sample 2(column 2) has blond hair. Legends show the distance of reconstructions from the original sample along the gradient steps, followed by overall classifier performance. Additionally, we fit a logistic function

$f(x) = \frac{1}{1+\exp^{-k(x-x_0)}}$ to each curve. All plots have been shift-aligned using x_0 . The confidence manifold for the same instance is fairly different across each model. As expected, the overall steepness increases as model trains to better discriminate samples. Intuitively, higher x_0 suggests that the classifier can easily discriminate the label with high confidence. For instance, for comparable overall accuracies, the manifolds suggest that model 2 has trained a decision boundary such that a manifold guided example is relatively close in image distance (compared to that of model 1). In the case of Sample 2, it is clear that model 2 mis-labels the data point with high confidence initially while learning to predict the correct label eventually. However, a decrease in x_0 as training progresses for both models suggests a significant shift of the decision boundary to be closer to Sample 2. Qualitative images corresponding to these manifolds are shown in the Appendix.

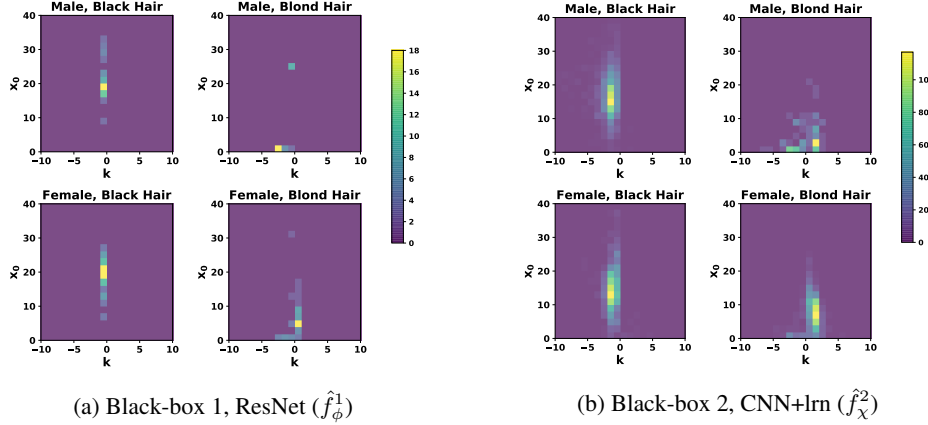


Figure 5: (a) and (b): 2d-Histograms of the parameters of the logistic function fits to the confidence manifolds for a ~ 4000 samples.

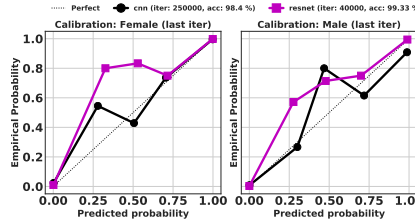


Figure 6: Reliability Diagram for Calibration stratified by (potentially protected) attributes of interest (gender): A perfectly calibrated classifier should manifest an identity function. Deviation from the identity function suggests mis-calibration and can be used for model comparison when accuracy and other metrics are comparable.

Figures 5a and 5b show the 2d histogram of the logistic function parameter estimates stratified by the target label and the attribute of interest (gender). This allows to summarize the confidence manifolds across groups of interest for overall model comparison. For reference, Figure 6 shows the Reliability Diagram for both black-boxes. The ResNet model generally demonstrates more uniform steepness across samples at different distances from the decision boundary compared to the CNN+lrn model. Both models have a relatively small x_0 for blond haired males suggesting lower confidence in their predictions. Thus, summarizing confidence manifolds provides additional insight that may not be characterized by Reliability Diagrams for model comparison.

6 Discussion

This work presents a novel approach to characterizing and explaining black-box supervised models via examples. An unsupervised implicit generative model is used as to approximate the data manifold, and subsequently used to guide the generation of increasingly confounding examples given a starting

point. These examples are used to probe the target black-box in several ways. In particular, we demonstrate the utility of manifold guided examples in automatically detecting bias in black-box learning w.r.t. a (potentially protected) attribute as well as for model comparison. The proposed method also allows one to visualize training progression and provides insights complementary to notions of calibration of the black-box model. Limitations of the proposed method include reliance on the implicit generator as a proxy of the data manifold. However, we note that we do not rely on specific architectures and/or training mechanisms for the generative model. We used images as they are easy to visualize even in high-dimensions. However extending our studies to complex datasets beyond images is a compelling future extension.

References

- [1] Agnar Aamodt and Enric Plaza. Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI communications*, 7(1):39–59, 1994.
- [2] Julius Adebayo, Justin Gilmer, Ian Goodfellow, and Been Kim. Local explanation methods for deep neural networks lack sensitivity to parameter values. *arXiv preprint*, 2018.
- [3] J Angwin, J Larson, S Mattu, and L Kirchner. Machine bias risk assessments in criminal sentencing. *ProPublica* <https://www.propublica.org>, 2016.
- [4] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- [5] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems*, pages 4349–4357, 2016.
- [6] Alison Callahan and Nigam H Shah. Chapter 19 - machine learning in healthcare. In Aziz Sheikh, Kathrin M. Cresswell, Adam Wright, and David W. Bates, editors, *Key Advances in Clinical Informatics*, pages 279 – 291. Academic Press, 2017.
- [7] R Dennis Cook and Sanford Weisberg. Characterizations of an empirical influence function for detecting influential cases in regression. *Technometrics*, 1980.
- [8] Morris H DeGroot and Stephen E Fienberg. The comparison and evaluation of forecasters. *The statistician*, pages 12–22, 1983.
- [9] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [10] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- [11] Been Doshi-Velez, Finale; Kim. Towards a rigorous science of interpretable machine learning. In *eprint arXiv:1702.08608*, 2017.
- [12] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.
- [13] Ethan Elenberg, Alexandros G Dimakis, Moran Feldman, and Amin Karbasi. Streaming weak submodularity: Interpreting neural networks on the fly. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017.
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [15] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. *arXiv preprint arXiv:1706.04599*, 2017.
- [16] Maya Gupta, Andrew Cotter, Jan Pfeifer, Konstantin Voevodski, Kevin Canini, Alexander Mangylov, Wojciech Moczydlowski, and Alexander Van Esbroeck. Monotonic calibrated interpolated look-up tables. *Journal of Machine Learning Research*, 2016.
- [17] Moritz Hardt, Eric Price, Nati Srebro, et al. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pages 3315–3323, 2016.

- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [19] Irina Higgins, Loic Matthey, Xavier Glorot, Arka Pal, Benigno Uria, Charles Blundell, Shakir Mohamed, and Alexander Lerchner. Early visual concept learning with unsupervised deep learning. *arXiv preprint arXiv:1606.05579*, 2016.
- [20] Michael C Hughes, Huseyin Melih Elibol, Thomas McCoy, Roy Perlis, and Finale Doshi-Velez. Supervised topic models for clinical interpretability. *arXiv preprint arXiv:1612.01678*, 2016.
- [21] Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*, 2015.
- [22] Been Kim, Rajiv Khanna, and Oluwasanmi O Koyejo. Examples are not enough, learn to criticize! criticism for interpretability. In *Advances in Neural Information Processing Systems*, pages 2280–2288, 2016.
- [23] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim. The (Un)reliability of saliency methods. *NIPS workshop on Explaining and Visualizing Deep Learning*, 2017.
- [24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [25] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. *arXiv preprint arXiv:1703.04730*, 2017.
- [26] Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. Visualizing and understanding neural models in nlp. *arXiv preprint arXiv:1506.01066*, 2015.
- [27] Zachary C Lipton. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*, 2016.
- [28] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [29] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *arXiv preprint arXiv:1706.07269*, 2017.
- [30] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
- [31] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. See <https://arxiv.org/abs/1610.02391> v3, 7(8), 2016.
- [32] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.
- [33] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [34] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [35] Pierre Stock and Moustapha Cisse. Convnets and imagenet beyond accuracy: Explanations, bias detection, adversarial examples and model criticism. *arXiv preprint arXiv:1711.11443*, 2017.
- [36] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *arXiv preprint arXiv:1703.01365*, 2017.

Appendix

xGEMs for MNIST

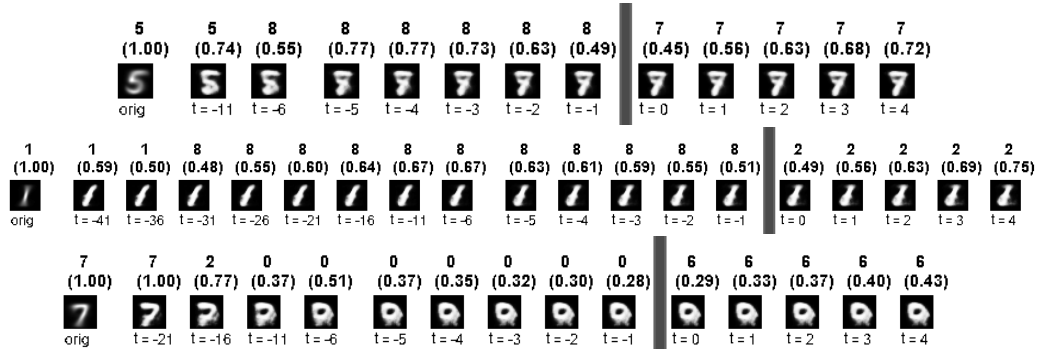


Figure 7: **xGEMs** for MNIST data. $\mathcal{G}_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{28 \times 28}$ is a VAE while the target black-box is a softmax classifier. Each row shows a manifold guided example transition for a single digit (labeled ‘orig’). The gray vertical bars indicate transition to the target label y_{tar} . Reconstructions in each row are intermediate reconstructions obtained using Algorithm 1. The confidence of the class prediction is shown in parentheses for each reconstruction.

Figure 7 shows manifold guided examples generated for a (multi-class) softmax classifier for MNIST⁷ digit data. The first row in Figure 7 shows manifold guided example for digit 5 if $y_{tar} = 7$, while second and third row show manifold guided examples for digits 1 and 7 with $y_{tar} = 2$ and $y_{tar} = 6$ respectively. Notice how while traversing the manifold, the classifier switches decision from 5 to 8 and then to the target label 7 (row 1). While the intermediate samples look like 7 to human eye, the classifier is biased toward predicting 8. Row 2 suggests a bias toward predicting 1 as 8 for a minor smudging (visible to human eye). Finally, the third row demonstrates how the manifold guided example for 7 suggests that the classifier considers a 0 to be labeled as 6. Thus manifold guided examples can provide insight into the decision boundary of the classifier for each pair of digits.



Figure 8: Training progression for celebA face image for the CNN+lrn model.

Case Study: Evaluating Model Training Progression

Figures 8 and 9 show **xGEMs** for the face corresponding to Sample 1 in Figure 4 for models CNN+lrn and ResNet respectively. Notice significant differences in the **xGEMs** and their trajectories even at comparable overall performance.

⁷<http://yann.lecun.com/exdb/mnist/>



Figure 9: Training progression for celebA face image for the ResNet model.