

ProtoShotXAI: Using Prototypical Few-Shot Architecture for Explainable AI

Samuel Hess

SHESS@EMAIL.ARIZONA.EDU

Gregory Ditzler

DITZLER@ARIZONA.EDU

Department of Electrical and Computer Engineering

University of Arizona

Tucson, AZ 85721, USA

Editor:

Abstract

Unexplainable black-box models create scenarios where anomalies cause deleterious responses, thus creating unacceptable risks. These risks have motivated the field of eXplainable Artificial Intelligence (XAI) to improve trust by evaluating local interpretability in black-box neural networks. Unfortunately, the ground truth is unavailable for the model’s decision, so evaluation is limited to qualitative assessment. Further, interpretability may lead to inaccurate conclusions about the model or a false sense of trust. We propose to improve XAI from the vantage point of the user’s trust by exploring a black-box model’s latent feature space. We present an approach, ProtoShotXAI, that uses a Prototypical few-shot network to explore the contrastive manifold between nonlinear features of different classes. A user explores the manifold by perturbing the input features of a query sample and recording the response for a subset of exemplars from any class. Our approach is the first locally interpretable XAI model that can be extended to, and demonstrated on, few-shot networks. We compare ProtoShotXAI to the state-of-the-art XAI approaches on MNIST, Omniglot, and ImageNet to demonstrate, both quantitatively and qualitatively, that ProtoShotXAI provides more flexibility for model exploration. Finally, ProtoShotXAI also demonstrates novel explainability and detectability on adversarial samples.

Keywords: Explainable AI, XAI, Local Interpretability, Few-Shot Learning.

1. Introduction

Large-scale neural networks have made promising strides to reduce the classification error over large volumes of data (He et al., 2016; Silver et al., 2016; Van Den Oord et al., 2016). For example, residual networks have surpassed human-level accuracy on ImageNet (He et al., 2016; Russakovsky et al., 2015a). Policy networks have defeated professional humans in the game of Go (Silver et al., 2016), and WaveNet can automatically generate natural human-sounding speech from text (Van Den Oord et al., 2016). Unfortunately, these black-box models often lack human interpretable reasoning for their inferences because of the millions of weights, biases, and nonlinear activation functions in their architectural complexity. One example in recent literature presented evidence that the top ImageNet classification models were biased on textures rather than the shape bias, demonstrating that prior explanations and interpretations of these model might be incorrect (Geirhos et al., 2019). Although textures may be sufficient for some tasks, the unexpected exploitation of texture leaves the

model vulnerable to undesirable predictions, especially with potential exposure to real-world data anomalies (Alcorn et al., 2019).

The concern of deleterious predictions from black-box models has triggered research efforts in eXplainable AI (XAI) that can be broadly organized into: (1) the development of explainable models that are comparable in prediction performance to their black-box counterparts, and (2) the development of techniques to derive explanations from existing black-box models (Bodria et al., 2021; Adadi and Berrada, 2018). Due to additional constraints of human interpretability, explainable models often struggle to achieve comparable performance with respect to their black-box counterparts. Additionally, explainable models are designed for a specific task, which limits their ability to extend to other tasks or datasets. In contrast, techniques that derive explanations from black-box models are support tools that can provide expectations on prediction behavior. The scope of explanations are often limited under the assumption(s) of the model and are not easily human interpretable. For example, some XAI approaches explain black-box models by weighting the importance of pixels. Unfortunately, these pixel attribution weights are still ambiguous because they require a human interpreter to extrapolate higher-level meanings from the weights (e.g., large pixel weights around a cat’s ears can indicate the shape of the ears are important high-level features). This human extrapolation of meaning is one of the core challenges with explaining black-box models: *there is no ground truth for the explanation itself*. In the cat example, large pixel weights around the ears could indicate that the shape of the ear is important for classification; however, the pixel weights could also indicate that the texture of the ear is what is truly important. Unfortunately, both of these explanations are just convenient for a human interpreter, but not necessarily correct. There could be many underlying true reasons for the pixel weights around the cat’s ears. For example, it could be the ratio of power levels between pixels, the color, or many other possibilities. The actual rationale for the pixel weight is unknown and a challenge in XAI.

A common approach to explain a model when ground truth does not exist is to involve actual human assessment through a survey. For example, Jeyakumar et al. (2020) performed Mechanical Turk surveys to determine user preference for explanations of deep neural networks (DNN) for image, text, audio, and sensor classification. Their study evaluated various state-of-the-art XAI against their proposed explain-by-example approach (namely, ExMatchina) and showed that users preferred ExMatchina for the majority of tasks (image, audio, and sensor). Although the study provided an interesting insight into human preference, there is a disconnect between a human’s explanation preference and the decision of the black-box model. That is, the empirical evaluation of user preferences can exploit the user’s confirmation bias (i.e., what the user already believes are important features).

Since there is no ground truth for the explanation of a black-box model’s prediction and human assessment of feature importance is inherently biased, we propose in this work to refocus XAI towards tractable tools for model exploration and trust. That is, instead of attempting to define a quantitative evaluation of a given explanation, we have developed a novel algorithm that is more amenable to exploring a model’s manifold space w.r.t input feature manipulations. Specifically, we present ProtoShotXAI, a flexible black-box exploration approach that takes advantage of a model’s trained feature representation layer, the weights of the classification layer (if available), and a Prototypical few-shot architecture to cross-compare a given query sample to class prototypes. Once the architecture is chosen,

a user can gain insight into the original model behavior and explore changes to the feature representation manifold by manipulating input features.

The primary contributions in this work can be summarized as:

- A novel application of few-shot architecture, namely ProtoShotXAI, was developed as a model-agnostic XAI approach that merges local interpretability methods and explain-by-example prototype methods. ProtoShotXAI leverages a novel cross-comparison of feature representation that allows for better analysis of a model’s decision.
- ProtoShotXAI is model-free in the sense that our approach is not specific to a particular neural network. Hence, our approach has broader applicability than some of the existing works in XAI.
- This work presents one of the first approaches to bridge XAI and few-shot architectures to increase the interpretability of few-shot neural networks and conventional models.
- Development of a tractable approach for model exploration and trust for black-box models, providing more insightful interpretations of real-world datasets.
- Publicly available code & data at <https://github.com/samuelhess/ProtoShotXAI/> which includes cross-comparisons between common XAI approaches with consistent representations as well as four unique experiments for XAI: adversarial MNIST, revolving six, Omniglot XAI, and a multi-class ImageNet example.

The remainder of this manuscript is organized as follows. Section 2 provides a background of recent and related works on XAI and few-shot learning. Section 3 describes the ProtoShotXAI approach, optimization, and use for learning localized regions of importance in images that explain a model’s decision from a feature representation. Section 4 presents experiments on three benchmark datasets against several state-of-the-art XAI approaches. Section 5 discusses the findings from our experiments and variants of ProtoShotXAI. Finally, Section 6 highlights the contributions, broader impacts, summary of the results and avenues for future research.

2. Related Work

2.1 eXplainable AI (XAI)

XAI can be broadly partitioned into two main categories: (1) approaches that are inherently explainable, and (2) approaches that derive explanations from black-box models (Bodria et al., 2021; Adadi and Berrada, 2018). In this work, we focus exclusively on explanations of existing black-box models that can be applied to DNNs. Within these methods, related work falls into the subcategory of locally interpretable model agnostic approaches. Popular approaches in this subcategory include saliency maps (Simonyan et al., 2013), Grad-CAM++ (Chattopadhyay et al., 2018), LIME (Ribeiro et al., 2016), SHAP (Lundberg and Lee, 2017), and Anchors (Ribeiro et al., 2018).

Saliency maps are visualized as heat maps that indicate feature importance and the maps are used to determine correlated areas in an image that are valuable to a model’s prediction.

Saliency maps are common amongst all the approaches mentioned above and originally evolved from literature in gradient-based methods. Similar to how a neural network trains, gradient methods take the derivative of an input w.r.t. an output class activation. The derivative provides a direction of feature manipulations that, when applied to the input features, would result in a positive activation gain for being “more like” the respective class. More formally, given a trained black-box classifier that transforms a D -dimensional sample $\mathbf{x} \in \mathbb{R}^D$ into an estimated vector of C classification scores $\hat{\mathbf{y}} \in \mathbb{R}^C$, the gradient map $\mathbf{z} \in \mathbb{R}^D$ is simply the partial derivative of the loss w.r.t. the input features $\mathbf{z} = \frac{\partial \mathbf{y}_c}{\partial \mathbf{x}}$. Unfortunately, unregularized gradient-based changes to the input often exploit the neural network’s function by providing a unique, unrealistic, and/or unexpected combinations of noisy gradients rather than providing clear insight into an input feature’s value w.r.t the network’s prediction. As a result, many variants of gradient-based approaches apply spatial averaging and other regularization techniques to constrain the gradients and reduce their noise. This regularization and smoothing helps to provide consistency to the scale of a feature in a spatial region which, in turn, improves human interpretability (Selvaraju et al., 2017; Chattopadhyay et al., 2018).

In contrast to gradient-based methods, Local Interpretable Model-agnostic Explanations (LIME) adds small perturbations on a give sample’s input features and create a linear model between the perturbations and the output predictions (Ribeiro et al., 2016). Specifically, a sample \mathbf{x} is transformed into a dataset $X' = \{\mathbf{x}'_1, \dots, \mathbf{x}'_N\}$ where \mathbf{x}'_i are N perturbations of \mathbf{x} . The dataset is passed through the classifier to produce a dataset of input samples and classification score pairings $T = \{(\mathbf{x}'_1, \mathbf{y}'_1), \dots, (\mathbf{x}'_N, \mathbf{y}'_N)\}$. A linear model can then be defined as $y'_i(c) = \mathbf{z}(c) \cdot \mathbf{x}'_i$ where the linear coefficients \mathbf{z} represent the map of feature importance for each class c and can be learned using multinomial logistic regression. In contrast to allowing gradients to dictate the direction of changes to the input, LIME effectively tests local input-output regions around a given sample more systematically. LIME assumes that, whereas the total input-output space is often highly complex and nonlinear, the small perturbation region around a given sample under test can likely be approximated as linear. This assumption, along with the significant variance of coefficients \mathbf{z} w.r.t the choice of input perturbations X' , have spawned many variants of LIME (Shankaranarayana and Runje, 2019; Zafar and Khan, 2019; Peltola, 2018; ElShawi et al., 2019; Bramhall et al., 2020). Of particular noteworthiness, the original authors of LIME also created an alternative rules-based approach called Anchors. In Anchors, the features of an input sample are again perturbed; however, unlike LIME, they are perturbed intelligently with a multi-armed bandit to produce a set of “anchors” that consistently hold for an output decision regardless of other (non-anchor) feature values. Instead of the feature attribution weights produced by other approaches, Anchors produce a set of logical feature boundaries for a decision. Often contrasted with LIME and Anchors is Shapley Additive Explanations (SHAP), which is an approach that focuses on the Shapley values from cooperative game theory. Shapley values define the importance of an input feature as the likelihood that the outcome would be different if a model was trained without that feature. Naturally, training numerous models with different features removed is often intractable, so the majority of SHAP approaches use clever ways to approximate Shapley values with different underlying assumptions on the type of data or model (Sundararajan et al., 2017; Smilkov et al., 2017; Erion et al., 2021). Notably, Expected Gradients and DeepSHAP (also known respectively

as GradientExplainer and DeepExplainer in the SHAP code libraries) are the versions of SHAP specifically designed for DNNs that we evaluate in this work. Expected Gradients use a subset of N training samples for background reference $T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ and define the explainable map as

$$\mathbf{z} = \mathbb{E}_{\mathbf{x}' \sim T, \alpha \sim U(0,1)} \left[(\mathbf{x} - \mathbf{x}'_i) \frac{\partial \mathbf{f}(\mathbf{x}'_i + \alpha(\mathbf{x} - \mathbf{x}'_i))}{\partial \mathbf{x}} \right]$$

where \mathbf{x}'_i is sampled from a subset of training samples, α is sampled from a uniform distribution, and the partial derivative is the gradient at a feature layer \mathbf{f} of the network w.r.t the input. The function $\mathbf{f}(\mathbf{x}'_i + \alpha(\mathbf{x} - \mathbf{x}'_i))$ represents a random hybrid between the input features under test \mathbf{x} and a sampled background \mathbf{x}'_i . \mathbf{f} can be the output of any feature layer up until the classification layer; however, the feature map becomes spatially coarser as the convolutional layers becomes coarser (and deeper layers in a CNN are often coarser). This observation means that the selection of a layer towards the middle should be preferred over a layer towards the end of the network.

Another related subcategory to our work are prototype methods. XAI, through prototypes, attempts to provide the user with a different but relatable sample (i.e., a “prototype”) that is an archetypical representation. Many of these approaches compute a distance between the prototypes and other samples in either the feature space of the input training data or a high-level feature manifold of a trained black-box model (Bien and Tibshirani, 2011; Kim et al., 2016; Koh and Liang, 2017; Papernot and McDaniel, 2018; Gurumoorthy et al., 2019). The distance metric is used to find a representative class sample that is close to the query sample. For example, ExMatchina uses the final feature layer (i.e., the layer next to the classification layer) to compute the cosine similarity between the query sample and training samples to find the N closest matches within that multi-dimensional manifold. Formally, given a trained black-box model and a sample \mathbf{x} , the exemplary sample \mathbf{x}'_e that best explains \mathbf{x} is

$$\mathbf{x}'_e = \arg \max_{\mathbf{x}' \in T} \frac{\mathbf{f}^T(\mathbf{x})\mathbf{f}(\mathbf{x}')}{\|\mathbf{f}(\mathbf{x})\|_2^2 \cdot \|\mathbf{f}(\mathbf{x}')\|_2^2}$$

where the \mathbf{x}' is from the the set of N training samples $T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ and $\mathbf{f}(\cdot)$ is the K -dimensional feature representation of \mathbf{x} . Intuitively, if \mathbf{x} is in the training set ($\mathbf{x} \in T$) then \mathbf{x}'_e is itself (i.e., $\mathbf{x}'_e = \mathbf{x}$).

2.2 Few-Shot Learning

In addition to related work in XAI, ProtoShotXAI leverages network architectures developed in the few-shot learning literature. Few-shot is a subcategory of classification where the goal is to categorize (with arbitrarily high accuracy) a set of unlabeled query samples into specific classes after being presented with a set of labeled support samples that contain one (or a “few”) exemplary sample(s) from each class. In contrast to conventional classification, where many training samples are available for each class, few-shot training datasets have a limited number of training samples per class. For example, the MNIST dataset has roughly 6,000 samples per class. In contrast, the comparable few-shot character dataset (known as Omniglot) has only 20 samples per class (Section 4 presents further details of these datasets).

Additionally, a benchmark evaluation of few-shot networks is performed on classes that are mutually exclusive from the training classes. Note that this is in contrast to conventional classification tasks that are only evaluated on classes in the training set. A few-shot task requires that one (or a “few”) exemplary samples of classes are provided during testing and conventional classification provides no such supplementary information.

The concept of few-shot has been around for decades in machine learning literature (Fink, 2004; Fei-Fei et al., 2006); however, many recent neural network developments towards few-shot have emerged in response to the creation of benchmark datasets. In response to the aforementioned Omniglot dataset, Siamese networks (Koch et al., 2015) created an architecture of two identical networks (i.e., networks with tied weights) where two different samples are presented at the input of each network. During training, Siamese networks compute a Euclidean distance at the feature layer of the network before the final sigmoid activation. The final activation is zero if two samples are from the same class and one if they are different (similar to the signature verification approach presented in Bromley et al. (1994)). The network was evaluated by randomly subsampling one exemplary sample from 20 random Omniglot classes and classifying the remainder of the samples from the 20 classes (known as 1-shot, 20-way evaluation). It achieved 93.4% classification accuracy.

Subsequent neural network developments extended the capability from one-shot to few-shot, improved the classification accuracy, and expanded the application to additional datasets. Notably, the authors of Matching Networks developed a network training scheme called episodic training where subsets of the training data are repeatedly (and randomly) selected to mimic the few-shot evaluation (Vinyals et al., 2016). Similar to Siamese networks, Matching Networks uses tied weights and computes a distance metric at the feature layer. Each of the exemplary samples produce a distance w.r.t the sample under test and a neural attention mechanism selects the most likely class. They demonstrated state-of-the-art classification performance on the Omniglot benchmark as well as the more complex mini-ImageNet dataset introduced by Ravi and Larochelle (Ravi and Larochelle, 2017). Most metric-based approaches continue to use the episodic training scheme, including Prototypical Networks. The authors of Prototypical Networks demonstrated that they could achieve even better performance than Matching Networks with a simplified architecture (similar to Siamese Networks), a Euclidean distance metric at the feature layer, and a softmax over the distances (Snell et al., 2017). Many few-shot architectures have been developed since Prototypical Networks, including metric architectures (Ye and Guo, 2018; Scott et al., 2018), generative architectures (Wang et al., 2018; Antoniou et al., 2018), meta learning architectures (Yoon et al., 2018; Finn et al., 2018), or some combination of these approaches; however, Chen et al. showed many of the few-shot methods have comparable results when using the same neural network backbone (Chen et al., 2019).

3. The ProtoShotXAI Architecture

XAI ExMatchina uses a distance function at the feature layer which is similar to the distance metric used in few-shot algorithms such as Matching Networks and Prototypical Networks. The objective for both ExMatchina and the few-shot approaches is similar as well: to find exemplary sample(s) that are close in the network’s feature manifold w.r.t the sample under test. The difference is that the few-shot architectures are designed to train a network for

contrastive classification and ExMatchina is designed to explain an existing trained network. Specifically, few-shot trains a network to minimize (at the feature layer) the distance between a query sample and an exemplary prototype whereas ExMatchina simply finds (at the feature layer of a trained classification network) the nearest training sample w.r.t the query. ExMatchina does not use an exemplary prototype nor does it use information from the trained classification weights.

Our ProtoShotXAI method is an XAI architecture that uses both an exemplary prototype and the information from the classification weights (when available) to explain a trained model’s decision. Specifically, ProtoShotXAI augments the feature layer with classification weights and computes a similarity score between an exemplary prototype and a query sample that needs explanation. Rather than finding the closest exemplary sample (as performed by ExMatchina), ProtoShotXAI computes an average exemplary prototype that has been shown to work well in few-shot learning research (Snell et al., 2017; Vinyals et al., 2016). When used in few-shot explanations, ProtoShotXAI computes the average prototype at the feature layer; however, when used with conventional classification networks, our ProtoShotXAI approach combines the feature layer with the weights of the classification layer (i.e., just prior to feature summation and activation applied at the last layer) which we demonstrate improves explainability over state-of-the-art XAI approaches (see Section 4). We argue that, whereas methods like ExMatchina and Expected Gradients remove layers of the classification network, ProtoShotXAI uses all the layers and is, therefore, more representative of the trained model and more applicable across different architectures.

After adding the few-shot architecture to a trained network, the similarity score can be used to provide explainable exemplary samples similar to ExMatchina. Additionally, the change in similarity score w.r.t. perturbations at the input provide insight into input feature importance and a feature attribution map can be produced (e.g., similar to the maps produced by LIME and SHAP). Unlike previous approaches, ProtoShotXAI provides the user with more flexibility to explore input feature changes and perturbations. For example, by using ProtoShotXAI on digits, the user can rotate a nine and observe the change in similarity score w.r.t classes nine and six (see Section 4). Further, unlike LIME (which also uses input perturbations), there is no separate model. The feature maps can be generated directly from the trained classification network.

The entire process to construct the ProtoShotXAI architecture and use it to create feature attribution maps is shown in Figure 1. The ProtoShotXAI process starts with creating a trained classification network or using a preexisting trained network. Section 3.1 formally describes the process of constructing the ProtoShotXAI architecture from the baseline trained model and Section 3.2 describes one specific process for perturbing the data to produce feature attribution maps using ProtoShotXAI. The feature attribution maps (which can be directly compared against other approaches), as well as additional uses for ProtoShotXAI, are demonstrated in Section 4.

3.1 ProtoShotXAI Modification of a Trained Network

ProtoShotXAI modifies a pretrained classification network into the dual network architecture commonly used by many few-shot approaches such as Siamese networks, Matching networks, and Prototypical networks. That is, the trained DNN is replicated up to the fea-

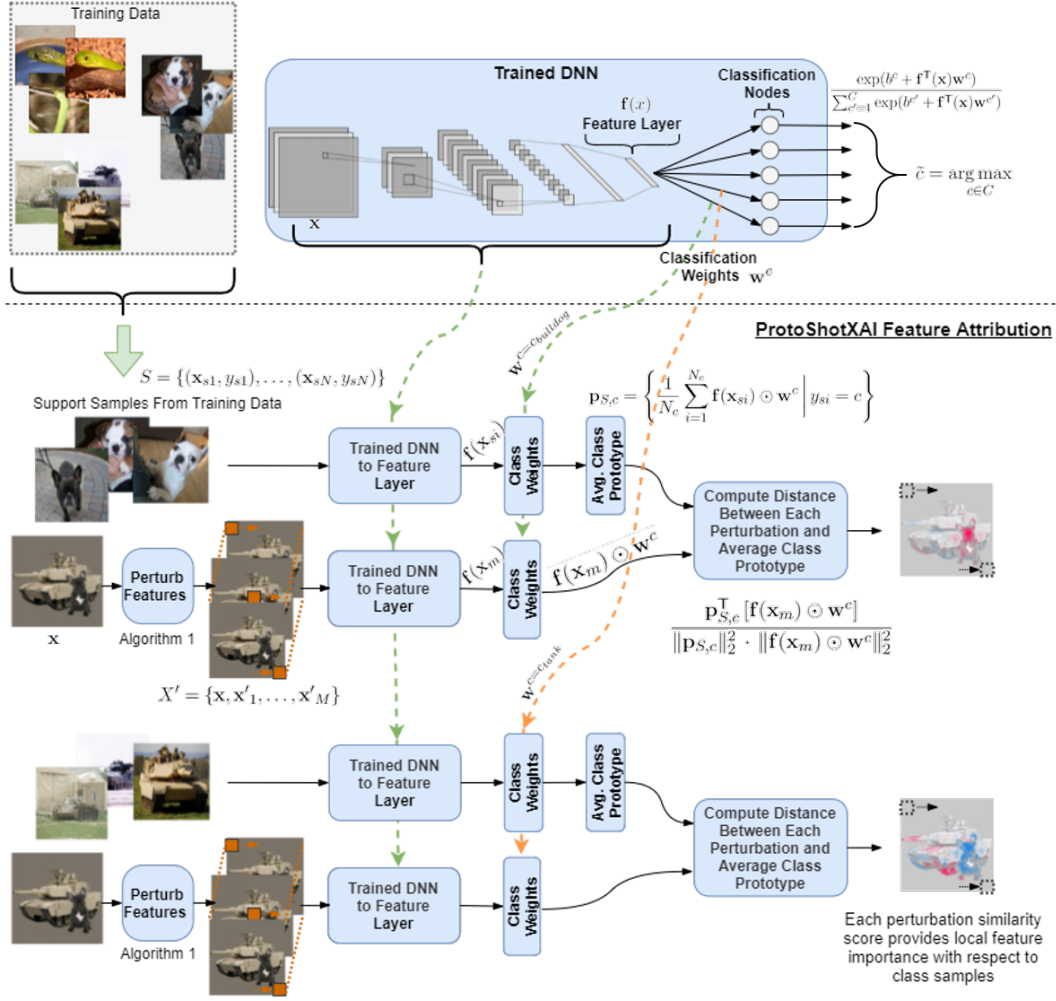


Figure 1: Overview of the architecture and feature attribution process of ProtoShotXAI. At the top of the figure is a trained DNN model and its corresponding training data. The ProtoShotXAI architecture has two branches. The first branch shows a set of support samples from the french bulldog passed through the trained DNN up to the the feature layer and multiplied by class weights of the french bulldog. A prototype vector is made from the average weighted features of the support samples. Similarly, the second branch contains the perturbations from the query sample (e.g., image containing the tank and french bulldog). For each pixel perturbation image, a cosine distance is computed between the prototype vector and the weighted image features and the output is the feature attribution score at each pixel. The process done for the french bulldog support set is repeated for the tank support set to give positive and negative attribution scores w.r.t each class.

ture layer for two branches: one branch contains the set of exemplary prototypes and the other branch contains variations of the query sample under observation. Formally, a classifier can be partitioned into a feature layer (the layer right before the classification layer) and a classification layer. The feature layer is represented by the nonlinear embedding function $\mathbf{f}(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^K$ which transforms a D -dimensional input sample $\mathbf{x} \in \mathbb{R}^D$ (such as an image) into a K -dimensional feature representation. At the classification layer, the features are multiplied by the trained static weights for each class node \mathbf{w}^c where $c \in \{0, \dots, C\}$ is a specific class out of the total number of class nodes C . The summation of the weighted features and bias b^c go through a softmax activation function for each of the classes and the estimated class \tilde{c} is the argmax of the activation as given by the following

$$\tilde{c} = \arg \max_{c \in [C]} \frac{\exp(b^c + \mathbf{f}^\top(\mathbf{x})\mathbf{w}^c)}{\sum_{c'=1}^C \exp(b^{c'} + \mathbf{f}^\top(\mathbf{x})\mathbf{w}^{c'})} \quad (1)$$

where $b^c \in \mathbb{R}$ is a scalar that represents the bias. As shown in the bottom of Figure 1, the first branch of the ProtoShotXAI architecture discards the activation, but keeps the trained network up until the feature layer and the weights of class. A subset $S = \{(\mathbf{x}_{s1}, y_{s1}), \dots, (\mathbf{x}_{sN}, y_{sN})\}$ of support samples from a class are from the training dataset Tr (e.g., images the french bulldog class). This support subset can be randomly selected (as done throughout this work) or, if the user knows a set of representative exemplary samples, they can be hand picked. All the N D -dimensional support samples represented by $\mathbf{x}_{si} \in \mathbb{R}^D$ are passed through the DNN up until the feature layer and are multiplied by the weights of the specific class of interest (e.g. the trained weights from the french bulldog class). The prototype feature vector of the support set $\mathbf{p}_{S,c}$ is the average feature multiplied by the weights for class c .

$$\mathbf{p}_{S,c} = \left\{ \frac{1}{N_c} \sum_{i=1}^{N_c} \mathbf{f}(\mathbf{x}_{si}) \odot \mathbf{w}^c \mid y_{si} = c \right\} \quad (2)$$

where \odot is the Hadamard product (i.e., an element by element multiplication operation). For the second branch of ProtoShotXAI, the sample under test \mathbf{x} (represented as the image containing both the tank and french bulldog in Figure 1) goes through a set of perturbations to observe the similarity response with respect to the prototype. Specifically, \mathbf{x} is transformed into a perturbation dataset $X' = \{\mathbf{x}, \mathbf{x}'_1, \dots, \mathbf{x}'_M\}$. This perturbed dataset can easily be created using a small changes to a different pixel in each sample. Similar to the support samples, each \mathbf{x}'_m goes through the DNN until the feature layer and is then multiplied by the weights of the specific class; however, no average is taken for the perturbations. Instead, each perturbation is compared against the prototype using a cosine similarity metric as represented in the following equation:

$$z_m(c, \mathbf{x}'_m) = \frac{\mathbf{p}_{S,c}^\top [\mathbf{f}(\mathbf{x}'_m) \odot \mathbf{w}^c]}{\|\mathbf{p}_{S,c}\|_2^2 \cdot \|\mathbf{f}(\mathbf{x}'_m) \odot \mathbf{w}^c\|_2^2} \quad (3)$$

ProtoShotXAI differs from ExMatchina in three distinct ways which are key contributions of our work for XAI: (1) the ProtoShotXAI distance function uses a prototype of features for it's reference rather than the samples from a single feature, (2) ProtoShotXAI multiplies the features by the weights of a specific class node to get more representative

features from the model, and (3) ProtoShotXAI uses input perturbations to cross-compare model responses to a baseline support set (like the feature attribution maps discussed in Section 3.2). In the case of few-shot approaches that already have a dual network architecture, the class weights can be set to a vector of ones and the same XAI approach can be used, making ProtoShotXAI more consistent across architectures.

3.2 Using ProtoShotXAI to Produce Feature Attribution Maps

After the construction of ProtoShotXAI, the entire architecture can be used in many ways to explore the feature space of the classification network. The ProtoShotXAI architecture receives two sets of inputs. One set are the support input features to compare against (e.g., samples from a single class) and the other can be any set of feature alterations to a query sample under test, such as adding perturbations, skewing features, adding noise, and/or setting features to zero. The change in the similarity score provides a measure for the effects of the feature alterations w.r.t to the class. That is, changing features that are more important to the network model will have a larger effect on the change to the similarity score.

As an example, input feature attribution maps can be created similar to the ones produced by LIME, SHAP, and Grad-CAM for images. We illustrate this implementation in the bottom of Figure 1 and present it as pseudo-code in Algorithm 1. The algorithm starts by choosing a set of input support image(s) (e.g., images of the french bulldog in Figure 1) and image under test (e.g., image containing both the tank and french bulldog in Figure 1). A support set can be created from random sample(s) of a class in the training data or any representative set of exemplary sample(s) and the image under test is any image to be evaluated w.r.t. the support. After the support and query data are selected, the algorithm calculates a baseline similarity metric $z_{ref}(c, \mathbf{x})$ for reference. The algorithm then loops through each pixel in the image and replaces the RGB value of that pixel (or small region surrounding that pixel) with an average value and computes the similarity metric for each spatial perturbation. Differences between the reference similarity score and the similarity score when the pixel is replaced is the feature attribution of that pixel. As the similarity score decreases, the pixel is weighted as being more important, and vice versa, as the similarity score increases.

4. Experiments

This section empirically demonstrates the proposed ProtoShotXAI approach on three benchmark datasets and against several state-of-the-art XAI approaches. We begin with a summary of the details related to implementing the datasets and the respective classification networks that ProtoShotXAI aims to explain. We then describe how ProtoShotXAI can be used to explore the original classification network qualitatively and quantitatively, and how our approach achieves comparable and often superior results to more complex algorithms for XAI. Throughout this section we use **boldface** text when referring to characters.

Algorithm 1 Implementation of ProtoShotXAI to Create Feature Attribution Maps

Input: Any representative set of N support image(s) $S = \{(\mathbf{x}_{s1}), \dots, (\mathbf{x}_{sN})\}$ for a single class c , a image under test \mathbf{x} , average RGB pixel vector \mathbf{R}_{RGB} , a trained classification model with frozen parameters

Output: Feature attribution map \mathbf{z} .

```
1:  $z_{ref} = \frac{\mathbf{p}_{S,c}^T [\mathbf{f}(\mathbf{x}) \odot \mathbf{w}^c]}{\|\mathbf{p}_{S,c}\|_2^2 \cdot \|\mathbf{f}(\mathbf{x}) \odot \mathbf{w}^c\|_2^2}$  // Reference score of  $\mathbf{x}$  w.r.t  $S$ 
2: for  $i = 1, \dots, I$  do // Loop through columns of the image
3:   for  $j = 1, \dots, J$  do // Loop through rows of the image
4:      $m = j + (i - 1)J$  // create a perturbation index for each 2D pixel
5:      $\mathbf{x}'_m = \mathbf{x}$  //  $\mathbf{x}'_m$  is initialized as  $\mathbf{x}$  before each perturbation
6:      $\mathbf{x}'_m(i, j) = \mathbf{R}_{RGB}$  // perturb the  $i, j$ th pixel to average RGB value
7:      $z(i, j) = z_{ref} - \frac{\mathbf{p}_{S,c}^T [\mathbf{f}(\mathbf{x}'_m) \odot \mathbf{w}^c]}{\|\mathbf{p}_{S,c}\|_2^2 \cdot \|\mathbf{f}(\mathbf{x}'_m) \odot \mathbf{w}^c\|_2^2}$  // Score for  $\mathbf{x}'_m$  w.r.t  $S$ 
8:   end for
9: end for
```

4.1 Implementation Details

Datasets In this work, we use three commonly used datasets to evaluate the interpretability of our approach: MNIST, Omniglot, and ImageNet. The MNIST dataset is a well-established benchmark for image classification. Although high classification performance can be achieved with very simple models, it is convenient to revisit this dataset in the context of XAI due to the overlap between features of different classes. For example, the number **eight** has a unique cross over “x”-like feature in the center; however, if the left half of the **eight** were removed then the remaining features appear like a **three**. Similarly, if the diagonal line from the bottom left to the top right of the **eight** were removed then the remaining features appear more like a **five**. We would expect similar conclusions when comparing XAI approaches on sufficiently trained classification networks.

The Omniglot dataset is a character database that is similar to MNIST but is used to benchmark few-shot algorithms. The dataset is composed of 50 different alphabets and a total of 1,623 character sets (Lake et al., 2011). For each of the 1,623 characters, there are only 20 human drawn samples which results in a database size of 32,460 samples. The data are originally 105x105 binary images and are resized to 28x28 to reduce input dimensionality. Additionally, the same 1,200 characters are rotated in multiples of 90 degrees to produce the set of 4,800 classes (96,000 samples) for training/validation and the remaining 1,692 classes (423 with 90° rotations) for testing (33,840 samples). To our knowledge the Omniglot dataset has not been used in any other XAI approach and its application is unique to our work.

ImageNet is a low-resolution image database that is also known as the ILSVRC-12 database (Russakovsky et al., 2015b). Samples vary in size but are often reduced to an image of size 3x224x224 RGB (as in this work) or less. The entire dataset is made up of nearly 14 million samples but the scale invariant feature set used in computer vision problems (and in this work) contains 1,000 classes with 1,200 samples per class (i.e., 1.2 million samples total). There is also a subset of the scale invariant ImageNet, known as Stylized-ImageNet, that we use in this work. Stylized-ImageNet substitutes textures in the

images with the goal of creating a training dataset that is biased towards shape instead of texture (Geirhos et al., 2019). For example, the fur from a dog is replaced with textures of skin from an elephant or bricks from a house while maintaining the shape and resemblance of the dog.

Classification Architectures All the XAI approaches evaluated in this work are applied to trained classification architectures, so each dataset was paired with a suitable neural network. For the MNIST dataset, we used the convolutional network shown in Table 1 which has a total of 1.2 million trainable parameters. For the Omniglot dataset, we followed the configuration used in prototypical networks (Snell et al., 2017) and shown in Table 2 with 112K trainable parameters. Finally, for ImageNet, we used the VGG16, ResNet50, and Xception network available in the Keras applications library. VGG16 was chosen as a common benchmark between methods because most comparable XAI approaches can be easily applied. VGG16 has 71% classification accuracy and 138 million parameters. Conversely, Xception network has a relatively small number of parameters (22 million) and is one of the best performing networks available in Keras (79% classification accuracy, third only to NASNetLarge with 89M parameters and 83% performance and InceptionResNetV2 with 56M parameters and 80% performance). Due to the complexity and depth of the Xception network (i.e., 126 layers with multiple convolutional layers, separable convolutional layers, and skip layers), only ProtoShotXAI and LIME can be easily applied to this network, and the other methods fail due to limited computational resources. Additionally, ResNet50 was used to compare conventional ImageNet training and a shape based variant of ImageNet, Stylized ImageNet. The network architectures were equivalent for both ImageNet and Stylized ImageNet, only the weights are different. VGG16 is shown in Table 3. For the specific ResNet50 and Xception network details, we refer the reader to the figures in the original papers (He et al., 2016; Chollet, 2017) or the supplemental material section on Github¹.

Classification Training and Optimization To train the MNIST and Omniglot networks, we used the ADAM optimizer (Kingma and Ba, 2015). The learning rate for the MNIST and Omniglot network was initially set to 6×10^{-5} and 1×10^{-3} , respectively. A scheduler was used for the Omniglot experiments to decrease the learning rate by half every 100 epochs for a total of 1,000 epochs. We did not use a scheduler in the MNIST network because sufficient performance was achieved with a constant learning rate after 100 epochs. For the VGG16, ResNet50 and Xception ImageNet, we used the available pretrained weights with the exception of one experiment where we trained the ResNet50 model from scratch using Stylized-ImageNet data intended to remove texture bias. For the ResNet50 Stylized-ImageNet training, the SGD optimizer was used with momentum. The initial learning rate was set to 0.1 and decreases by 90% every 20 epochs for a total of 40 epochs.

XAI Experimental Protocol The experimental protocol was to train each network with their respective data to achieve sufficient performance. After the network was trained, we constructed the ProtoShotXAI architecture presented in Section 3.1 and demonstrated in Figure 1. We used the feature layer presented in bold for each network in each of the respective tables. When available, we applied the weights of the classification nodes to

1. <https://github.com/samuelhess/ProtoShotXAI/>

Model: Convolutional MNIST Model

Layer	Output Shape	# of Parameters
Input	(28, 28, 1)	0
Conv2D (ReLU Activation)	(26, 26, 32)	320
Conv2D (ReLU Activation)	(24, 24, 64)	18496
MaxPooling2D	(12, 12, 64)	0
Dropout 25%	(12, 12, 64)	0
Flatten	(9216)	0
Dense Feature Layer (ReLU Activation)*	(128)	1179776
Dropout 50%	(128)	0
Dense Classification Layer (Softmax Activation)	(10)	1290

Total parameters: 1,199,882

Trainable parameters: 1,199,882

Non-trainable parameters: 0

* The output of this layer corresponds to $\mathbf{f}(\mathbf{x})$ in Section 3.

Table 1: Convolutional MNIST Model

Model: Prototypical Omniglot Model

Layer	Output Shape	# of Parameters
Input	(28, 28, 1)	0
Conv2D (BatchNorm + ReLU Activation)	(28, 28, 64)	896
MaxPooling2D	(14, 14, 64)	0
Conv2D (BatchNorm + ReLU Activation)	(14, 14, 64)	37184
MaxPooling2D	(7, 7, 64)	0
Conv2D (BatchNorm + ReLU Activation)	(7, 7, 64)	37184
MaxPooling2D	(3, 3, 64)	0
Conv2D (BatchNorm + ReLU Activation)	(3, 3, 64)	37184
MaxPooling2D	(1, 1, 64)	0
Flatten (Feature Layer)*	(64)	0

Total params: 112,448

Trainable params: 111,936

Non-trainable params: 512

* The output of this layer corresponds to $\mathbf{f}(\mathbf{x})$ in Section 3.

Table 2: Prototypical Omniglot Model

Model: VGG16 ImageNet Model

Layer	Output Shape	# of Parameters
Input	(244, 244, 3)	0
Conv2D (BatchNorm + ReLU Activation)	(224, 224, 64)	1792
Conv2D (BatchNorm + ReLU Activation)	(224, 224, 64)	36928
MaxPooling2D	(112, 112, 64)	0
Conv2D (BatchNorm + ReLU Activation)	(112, 112, 128)	73856
Conv2D (BatchNorm + ReLU Activation)	(112, 112, 128)	147584
MaxPooling2D	(56, 56, 128)	0
Conv2D (BatchNorm + ReLU Activation)	(56, 56, 256)	295168
Conv2D (BatchNorm + ReLU Activation)	(56, 56, 256)	590080
Conv2D (BatchNorm + ReLU Activation)	(56, 56, 256)	590080
MaxPooling2D	(28, 28, 256)	0
Conv2D (BatchNorm + ReLU Activation)	(28, 28, 512)	1180160
Conv2D (BatchNorm + ReLU Activation)	(28, 28, 512)	2359808
Conv2D (BatchNorm + ReLU Activation)	(28, 28, 512)	2359808
MaxPooling2D	(14, 14, 512)	0
Conv2D (BatchNorm + ReLU Activation)	(14, 14, 512)	2359808
Conv2D (BatchNorm + ReLU Activation)	(14, 14, 512)	2359808
Conv2D (BatchNorm + ReLU Activation)	(14, 14, 512)	2359808
MaxPooling2D	(7, 7, 512)	0
Flatten	(25088)	0
Fully Connected Layer	(4096)	102764544
Fully Connected Feature Layer*	(4096)	16781312
Dense Classification Layer (Softmax Activation)	(1000)	4097000
Total params: 138,357,544		
Trainable params: 138,357,544		
Non-trainable params: 0		

* The output of this layer corresponds to $\mathbf{f}(\mathbf{x})$ in Section 3.

Table 3: VGG16 ImageNet Model

the features for the respective class. For comparison, each approach is used to explain the same image with the same model. Many approaches display feature attribution weights differently, but we used the approach developed by SHAP for all methods to maintain consistency for comparison. That is, once the feature attribution weights \mathbf{z} are computed, we compute the 99.9 percentile of the absolute values of \mathbf{z} . The negative and positive of this value represents the dynamic range of the color axis where red and blue colors indicate positive and negative feature attributions, respectively.

4.2 MNIST

In this section, we demonstrate three different aspects of ProtoShotXAI with experiments for the MNIST dataset. In the first experiment, we qualitatively evaluate the locally interpretable algorithms by visualizing the feature representations of the digit **eight**. **Eight** is a unique digit due to the subset of digits that are contained within it. For example, the digit **three** is commonly contained as a subset of **eight** as well as parts of **five**, **six**, and **nine**, making it convenient for cross-class explanations. In the second experiment, we demonstrate how ProtoShotXAI can be used to quantitatively assess the classification change as a **six** is rotated until it is classified as a **nine**. Finally, the last MNIST experiment analyzes the XAI approaches on an adversarial dataset. The adversarial dataset is unique because the true class and the predicted class from the classifier are often different and the features of ProtoShotXAI provide insight into the network’s classification errors.

4.2.1 FEATURE ATTRIBUTION ON EIGHT

In the first experiment, we qualitatively evaluated ProtoShotXAI with an MNIST sample of the digit **eight**. The digit **eight** was chosen as a test sample due to the presence of multiple digits contained within the shape of the **eight**. The same **eight** was tested against each of the digit classes (0-9) using the six state-of-the-art XAI algorithms, namely, ProtoShotXAI, Deep SHAP, Expected Gradients, Grad-CAM, LIME, and ExMatchina. The reference pixel value used in ProtoShotXAI is set to zero which indicates lack of digit presence in a pixel. The results of each algorithm were normalized using the same standard applied in the SHAP work, which is necessary for a fair comparison, and the results are displayed in Figure 2.

One of the first observations from Figure 2 is that pixels with negative feature attributions are much more prevalent in ProtoShotXAI over the other approaches. The SHAP variants are the only other methods with negative feature attributions; however, the magnitude of those attributions are relatively small. The prominent negative features of ProtoShotXAI are attributed to the weights of the classification nodes which cause strong negative responses to counterfactual features of the given class (see Section 4.2.3 for more examples with adversarial perturbations). As a result, the remaining positive features in ProtoShotXAI highlight expected areas of interest. For example, the positive features for **two** are symmetrical to **five** about a vertical axis and the positive features for **six** are symmetrical to **nine** about a horizontal axis. With respect to the digit **three**, there is significantly more clarity in the ProtoShotXAI positive and negative features when compared to other methods. Additionally, ProtoShotXAI is more spatially precise than other methods. SHAP and Grad-CAM use only part of the convolution network layer and rely on spatial interpolation which causes a spatial blurring effect. LIME uses spatial segmentation



Figure 2: Example of feature attribution weights of various approaches on an exemplary character **eight** (shown in gray throughout). Consistent with the SHAP conventions, red areas indicate positive features attributions and blue areas indicate negative feature attributions. ProtoShotXAI shows more contrast in features (e.g., positive attributes for digits **two**, **three**, **five**, and **seven** as sub-features of the digit **eight**) compared to other methods. ExMatchina is the one prototype method shown due to the mathematical similarity to ProtoShotXAI. ExMatchina only provides the nearest exemplary samples in a class which is apparent from provided exemplary digits like **three**, **five**, and **nine** but does not provide the same analytical clarity of feature attributions.

areas that cause entire regions to be flagged as a positive or negative feature attribute. ExMatchina is the only explain-by-example approach shown due to its mathematical similarity to ProtoShotXAI but is not intended to provide cross-class comparisons. Hence, the closest samples ExMatchina provide are arguably “**eight**-like” (see exemplar **three**, **five**, and **nine**) but do not provide feature clarity like the other approaches in this context.

4.2.2 REVOLVING SIX

We also performed a simple – yet meaningful – experiment that shows ProtoShotXAI can be used for more than feature attribution maps, the similarity score it produces can bring fidelity to classification. Further, ProtoShotXAI has better agreement than the class of the nearest exemplary sample (i.e., ExMatchina). To demonstrate these properties of ProtoShotXAI, we consider an image of the digit **six** that is rotated from 0° to 360° . If a **six** is rotated 180° about the center of the digit, the image can be incorrectly classified as a **nine** by the network. It is worth mentioning that most humans would also incorrectly classify the rotated **six** as a **nine**. A question that arises from this scenario is: as a **six** is rotated from 0° and 180° , when does this transition occur that causes a different classification, and how rapid is the transition?

In this experiment, we rotate a **six** counter clockwise from 0° to 359° in increments of 1° then pass each rotated sample to a DNN. The model classifying these different rotational views of the **six** only predict three different classes throughout the rotation: a **six**, a **zero**, and a **nine**. The transitions occur at 93° (transition from 6 to 0), 158° (transition from 0 to 9), 278° (transition from 9 to 0), and 329° (transition from 0 to 6). The top plot in Figure 3 shows the transition of the class predicted by the model (solid blue line), the ProtoShotXAI score (dashed red line), and the ExMatchina* score (dashed green line). For each of these rotated images, we use the ProtoShotXAI architecture to compute the distance between the rotated sample and a prototype of 100 random exemplary samples from each digit class of interest (i.e., classes **six**, **zero**, **five**, and **nine**). ProtoShotXAI scores for these classes of interest as the **six** is rotated are shown in the middle plot of Figure 3. For comparison, the same 100 random exemplary samples are used in ExMatchina to retrieve the best cosine similarity within each class, which are shown in the bottom plot of Figure 3. Note that this evaluation of ExMatchina was not proposed in the authors’ original implementation. Since we revised ExMatchina from an explain-by-example approach to an analytical approach closer to ProtoShotXAI, we represent it as ExMatchina*.

The ProtoShotXAI and ExMatchina* scores have similar behaviors that are correlated to the model’s predictions; however, these scores have better explainable fidelity. For example, the scores from ProtoShotXAI show a transition between classes of approximately 20° - 30° before the model is certain about the next class. The difference between ProtoShotXAI and ExMatchina* in this synthetic example is that ProtoShotXAI finds the cosine similarity between the rotated **six** and the *average prototype* of the 100 exemplary samples, whereas ExMatchina* computes the cosine similarity between the rotated **six** and all the 100 exemplary samples then chooses the *prototype with the maximum similarity*. The result from ProtoShotXAI is more stable because the reference class prototypes will not change as the **six** rotates; however, this observation can not be said about the response of ExMatchina*’s score. The stability of ProtoShotXAI is further demonstrated by comparing the class predicted by the model (shown in the top plot of Figure 3) with the class that corresponds to the maximum score for ProtoShotXAI (dashed red line) and the maximum cosine similarity used in Exmatchina* (dashed green line). The output of the predicted class from ProtoShotXAI scores is more consistent to the class predicted by the model whereas the ExMatchina* cosine distance produces more deviations and even temporarily predicts a **five**. Therefore, we conclude ProtoShotXAI produces a similarity score that brings explainable fidelity to a model’s prediction and provides more insight than ExMatchina*.

4.2.3 ADVERSARIAL MNIST

In the context of classification, adversarial data are crafted to exploit a model’s behavior to produce incorrect predictions by adding a small perturbation to the original image that is imperceivable to a human but fools the model (Goodfellow et al., 2014; Sadeghi et al., 2020; Qiu et al., 2021). For example, specific small perturbations to an MNIST digit can appear to be minor alterations to the human eye yet cause drastic changes to the output prediction of a trained neural network. This behavior causes hesitancy integrating neural networks in precarious situations where an adversary might be able to inflict malicious data with the intent of causing deleterious predictions. As a result, a relevant goal of XAI is to

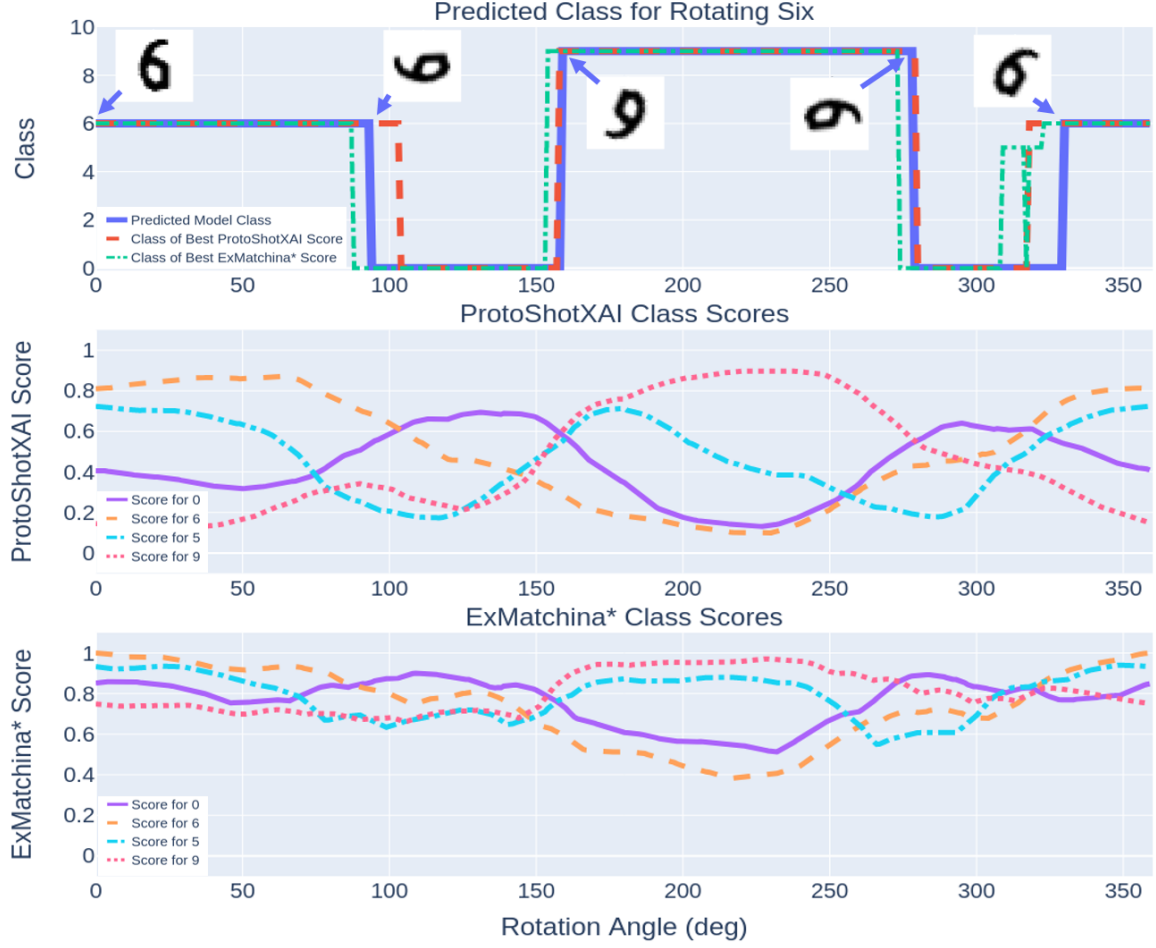


Figure 3: Predicted class as a **six** is rotated from 0° to 359° (blue line in the top plot) and corresponding ProtoShotXAI and EXMatchina* scores for classes **zero**, **six**, and **nine**. The red and green lines in the top plot result from assigning the class of the largest similarity score at each rotation angle for ProtoShotXAI and ExMatchina*, respectively. ProtoShotXAI has more consistency with the model predictions due to the prototype averaging.

explain when and why adversarial data corrupts model behavior. If XAI can explain when and why, then deleterious predictions from an adversary can be bounded or mitigated which increases user trust in the black-box model.

To explore the adversarial behavior of ProtoShotXAI, we crafted an adversarial dataset by taking the trained neural network from Table 1 and training it on the MNIST dataset. The adversarial samples are generated from the testing dataset using the *Fast Gradient Sign Method* (FSGM), which generates adversarial samples of the form (Goodfellow et al., 2014):

$$\mathbf{x}_{\text{adv}} = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} L(\theta, \mathbf{x}, y)) \quad (4)$$

where $\epsilon > 0$ is the adversarial budget, θ are the parameters of the network and $J(\cdot, \cdot, \cdot)$ is the categorical cross-entropy loss. The adversarial samples are presented to the ProtoShotXAI network.

The results of this experiment are shown in Figures 4 and 5. There are several observations that we make by considering the results in Figure 4. The feature attribution maps of all methods on the adversarial data were noisy and provided little-to-no insight; however, by plotting the component weights of the support prototype, the query sample, and the ProtoShotXAI score, we noticed a clear difference between the adversarial samples and benign samples. A benign **four** and adversarial **four** were selected from the testing dataset to demonstrate this difference. As shown in the top two plots of Figure 4, a benign **four** correlates well to the support prototype of a **four** (top left plot), which causes prominent peaks in the ProtoShotXAI components and subsequent scores. The prototype of a **five** causes the query **four** to produce a lot of negative weights and is not correlated well to the support prototype of a **five** (top right plot). The result is as expected: the benign **four**'s score for the class **four** prototype is much larger than the class **five** prototype. However, in the adversarial query sample of a **four**, neither the score for the prototype of the **four** or the prototype of the **five** are well correlated (bottom left and right plots, respectively). The overall score is poor for both classes but is marginally worse for class **four** than for class **five**. As a result, the adversarial sample is incorrectly classified as the **five** class.

The variation between the adversarial and benign scores for the in-class samples in Figure 4 is novel to the ProtoShotXAI approach and begs the question: *can ProtoShotXAI be useful to identify adversarial samples?* From the analysis of ProtoShotXAI scores, the results show that the adversarial samples corrupt the network at the feature layer by decorrelating the feature weights from the classification weights. The features in a benign sample strongly align to the weights of its own class, whereas the adversarial features do not strongly align with any class. This lack of alignment to a class causes minor correlations to decide the prediction. This result is made clear by running an experiment that generates the in-class ProtoShotXAI scores for 1,000 random samples of both the benign and adversarial MNIST data. The benign and adversarial score distributions are shown in the left plot of Figure 5. As expected, the adversarial digits have a significantly lower average than the average ProtoShotXAI scores on benign digits. It is also worth noting that the variance of ProtoShotXAI's score is larger than its benign counterpart. A simple detector can be created by applying a threshold to the ProtoShotXAI score. The right plot of Figure 5 shows the receiver operating characteristic (ROC) curve for adversarial detection versus falsely detecting benign digits as adversaries. The results of the simple threshold detector

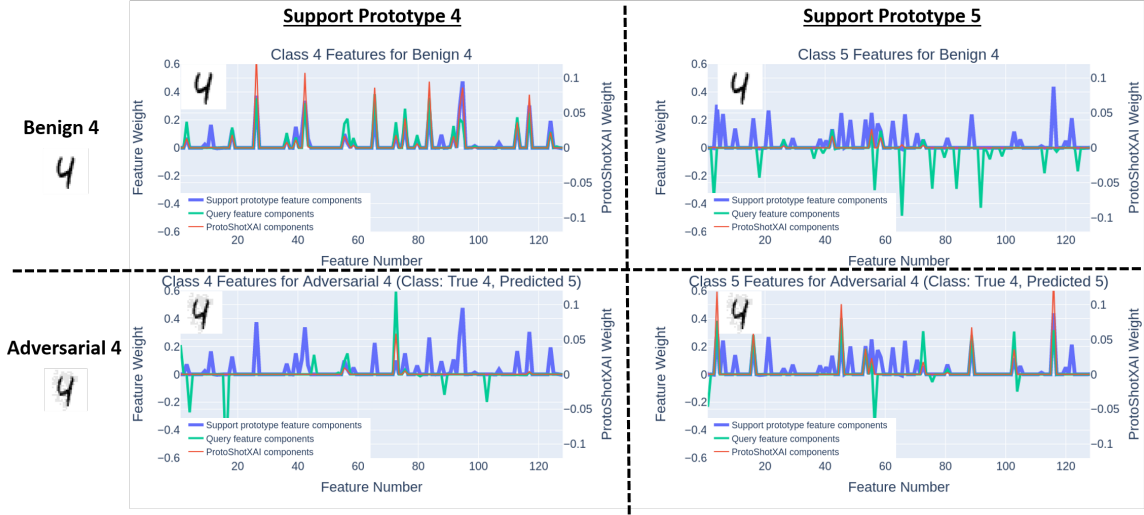


Figure 4: Results on the adversarial MNIST dataset which shows the decomposition of features from a benign **four** (top) and adversarial **four** (bottom) w.r.t. class **four** prototype (left) and class **five** prototype (right). From these plots, the benign **four** has a very positive response to the class **four** prototype and a negative response to the class **five** prototype. The adversarial samples do not show this same behavior, it is more uniformly distributed across features. ProtoShotXAI can demonstrate at this level that the features are corrupted, but what components caused feature corruption is still unanswered.

demonstrates that ProtoShotXAI can answer when a sample is adversarial and helps to improve model trust; however, it does not directly answer why the sample is adversarial (i.e., why from an interpretability perspective). We leave the analysis of “why” the sample is adversarial as future work and this question of “why” is an ongoing area of research (Merrer and Trédan, 2020; Tomsett et al., 2020).

4.3 Omniglot

ProtoShotXAI borrows from a few-shot architecture which makes it easy to apply to few-shot classification methods. To the best of our knowledge, ProtoShotXAI is the first method to explicitly apply XAI to datasets such as Omniglot. Since many few-shot architectures use the dual network structure for support and query samples, applying ProtoShotXAI to few-shot tasks is equivalent to a standard neural network with the classification weights set to a vector of ones. Figure 6 shows a character from the Omniglot database (left) that can be classified with the trained model from Table 2. Similar to MNIST, we used ProtoShotXAI with Algorithm 1 and a reference pixel value of zero which indicates lack of digit presence in a pixel. Although many pixels have a positive in-class response for the Omniglot character (right), it is observed that the order of importance appears to be the top left, the bottom, and the top right sub-character. Additionally, the specific size of the top right sub-character appears to be less important than the spatial center based on the strong positive response from sub-character’s spatial center.

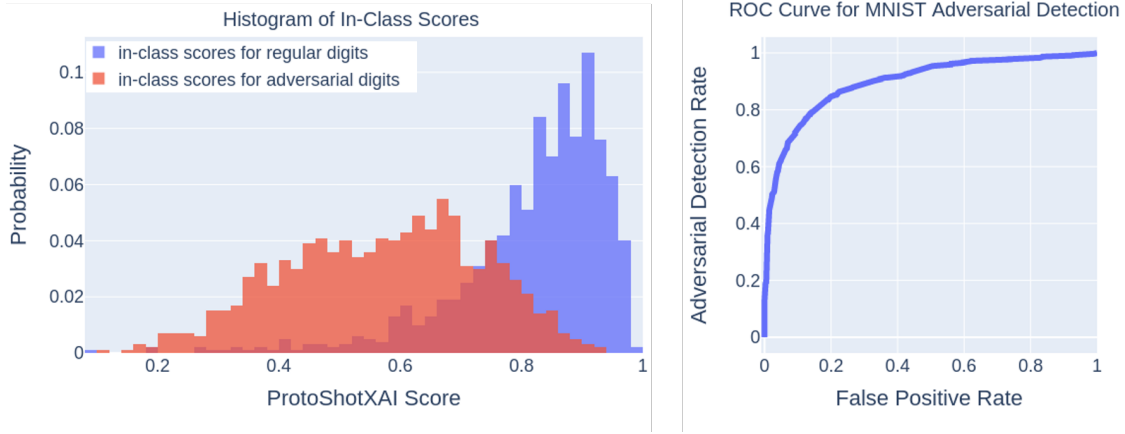


Figure 5: The plot on the left shows the histogram of the in-class scores for the benign **four** (in blue), the adversarial **four** (in orange), and the overlap (in red). The plot on the right is the ROC curve for a simple histogram threshold detector applied to the ProtoShotXAI scores. The result is the ProtoShotXAI scores can be used to detect adversarial samples by knowing the in-class scores from a validation dataset.

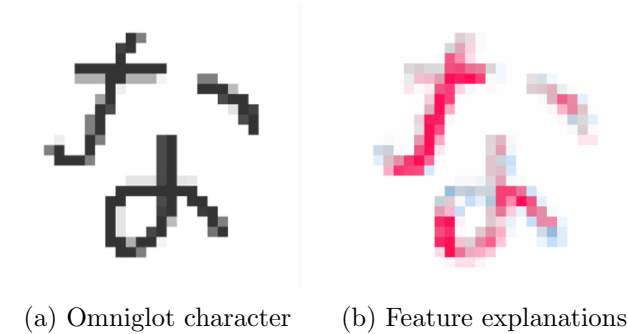


Figure 6: An example of an Omniglot character (left) and ProtoShotXAI feature attributions w.r.t in-class prototype using Algorithm 1 (right). The features indicate strongest weighting to the sub-character on the top left and weakest weighting to the sub-character on the top right. We demonstrate similar conclusion when looking at the similarity score w.r.t removing the larger sub-characters as shown in Figure 7.

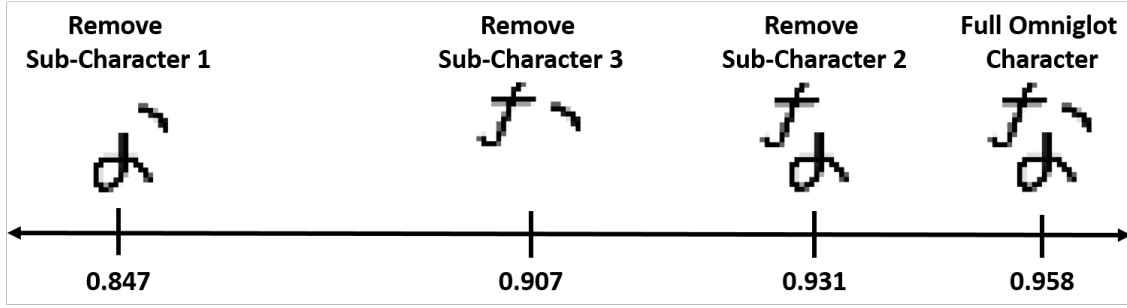


Figure 7: ProtoShotXAI similarity score w.r.t in-class prototype. The image on the far right with the highest similarity score corresponds to the full Omniglot character where the three images on the left correspond to the ProtoShotXAI score with different sub-characters removed. Removing the top left sub-character causes the largest negative deviation from the similarity score, indicating it is the most important sub-character of the three. The order of sub-character importance is in agreement with the feature attribution maps shown in Figure 6.

In contrast to feature attributions, ProtoShotXAI can evaluate the response to large scale changes to features. Using the previous character as an example, we removed each of the respective sub-characters to determine the overall network response to its class using the ProtoShotXAI score. Figure 7 shows the resultant scores on a line plot. The overall score dropped the most by removing the top left sub-character, indicating that was most important to the overall character classification. Conversely, removing the top right sub-character caused little change to the overall score, indicating that it was least important to the overall character classification. This result corresponds to the same order of sub-character importance obtained from the feature attribution maps.

4.4 ImageNet

The previous sections evaluated ProtoShotXAI against state-of-the-art digit classification methods (MNIST) and also demonstrated that ProtoShotXAI could be used for few-shot character classification (Omniglot). Both the MNIST and Omniglot datasets are images with approximately binary pixels; however, many real-world applications have complex features or pixels. The experiment in this section aims to qualitatively demonstrate that the ProtoShotXAI approach can be applied to datasets with complex image features, namely ImageNet. To control the experiment and demonstrate positive and negative feature attributions, we created an artificial input of an image containing a tank and a french bulldog for an explanation (i.e., two classes within the 1,000 ImageNet classes). The background is an average pixel value that appears as a brownish color. There are a total of four pretrained classification networks used in the assessment: VGG16, Xception Network, and ResNet50, all with ImageNets weights, and ResNet50 with Stylized-ImageNet weights. “ImageNet weights” refer to the default weights provided by Keras and the “Stylized-Imagenet weights” refer to weights from Stylized-Imagenet training as discussed in Section 4.1. Stylized-ImageNet substitutes textures in the images with the goal of training a network that is biased to-

wards shape instead of texture. Figure 8 shows the comparison of ProtoShotXAI, Expected Gradients (SHAP Gradient Explainer), Grad-CAM, and LIME for each trained network. Only LIME and ProtoShotXAI were applicable on Xception Network due to the network’s complexity. The gradient-based methods (i.e., Expected Gradients and Grad-CAM) operate at a specific convolutional layer and break down with unconventional convolutional architectures (like that of Xception Network). In contrast, LIME and ProtoShotXAI use the full model and do not have the same limitation. In this work, Grad-CAM operates at the last convolutional layer and Expected Gradients operates on layers 7 and 38 for the VGG16 and ResNet architectures, respectively. The choice of the convolutional layer is a hyperparameter in Expected Gradients where deeper layers tend to provide a better representation at the cost of lower spatial resolution. In this work, the layer was selected based on a comparable spatial resolution to the other methods, which results in a fair comparison.

For each trained network, ProtoShotXAI demonstrates consistent positive features (red) for the tank attributions around barrel/treads and negative features (blue) in the region of the french bulldog (see Figure 8). For the french bulldog attributions, ProtoShotXAI has positive responses around the face and negative responses scattered throughout the tank. All methods appear to have inconsistencies w.r.t attribution and class (e.g., some positive feature attributions on the tank for the french bulldog class); however, the response of ProtoShotXAI is more spatially consistent throughout. Expected Gradients is the closest to ProtoShotXAI for the VGG16 class with more “incorrect” positive features on the tank for the french bulldog class. Additionally, Expected Gradient’s dot-like appearance on ResNet50 is from the method’s use of deeper convolutional layers that partition the image into low-resolution spatial features. In contrast, the other gradient method, Grad-CAM, has a more blurred appearance and only identifies positive features of interest. The blurred appearance causes background regions to be positively indicated as feature attributions as well. LIME does not have the same issue of positively indicated feature attributions due to a segmentation algorithm that causes very defined region edges; however, LIME also segments large regions inconsistently, such as the positive feature attributions around the french bulldog’s feet w.r.t the tank class.

ResNet50 with Stylized-ImageNet is unique to this experiment because it is the only network trained with the modified ImageNet dataset meant to remove the network’s texture bias while maintaining the shape bias. We expected to see different network responses that indicate isolation of feature boundaries rather than regions, and unfortunately, this behavior was not explicitly observed. The feature attribution maps for ResNet50 trained with Stylized-ImageNet are drastically different from ResNet50 trained with the original ImageNet but it is inconclusive to argue that any XAI feature attribution methods show evidence of a shape bias. Although considered out of scope for this work, our XAI results warrant further investigation to determine if ResNet50 (trained with Stylized-ImageNet) has learned some complex features other than explicit shapes, such as the underlying textures from the stylized transformation.

5. Discussion

Throughout the experiments section, ProtoShotXAI was demonstrated to be a flexible model exploration tool that produces comparable (and often superior) qualitative and quan-

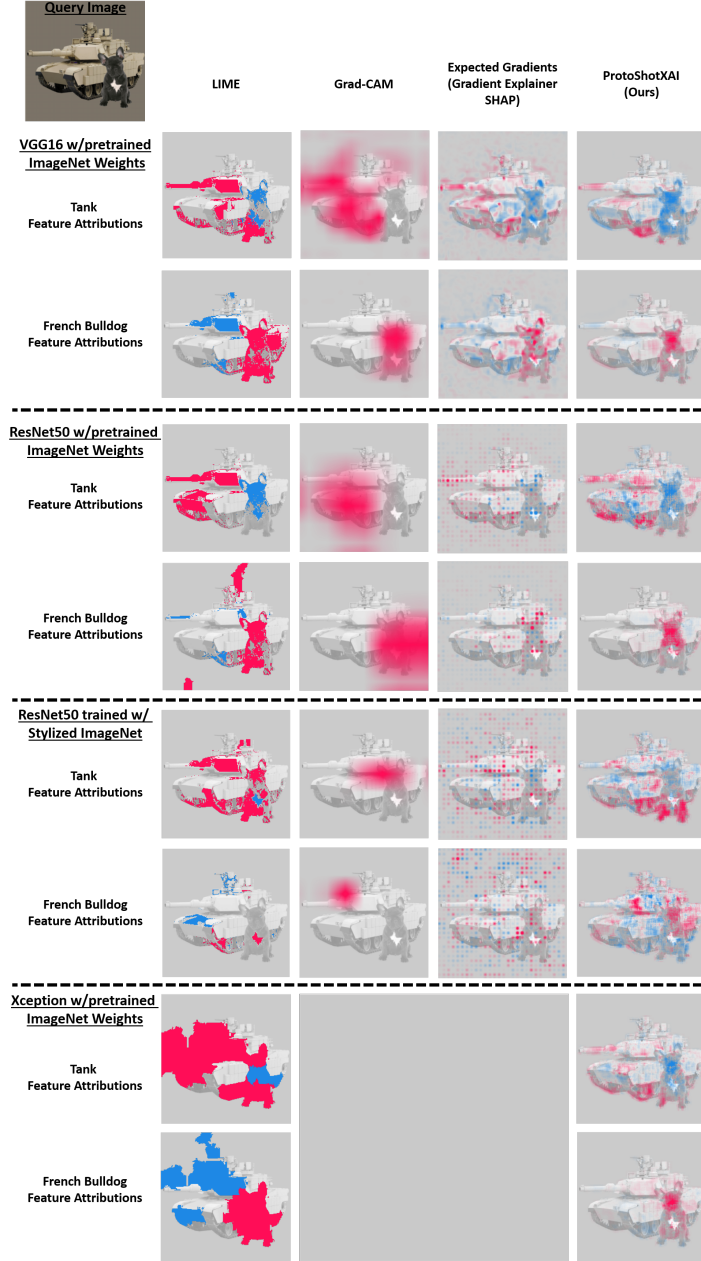


Figure 8: Comparison of feature attributions produced by four different methods (ProtoShotXAI, Expected Gradients SHAP, Grad-CAM, and LIME) and four different networks. The query image containing the tank and french bulldog is shown in the top left. Each row corresponds to a feature attribution of a given network for a specific class (i.e., tank or french bulldog class). Each column corresponds to a particular approach. Due to the complexity of the Xception Network architecture, the gradient methods, which expect a specific convolutional architecture, were not applicable. The feature attributions from ProtoShotXAI are noticeably more consistent (i.e., positive and negative feature attributions are more isolated to their respective class region).

titative results. ProtoShotXAI is uniquely straightforward to implement since it does not require another trained model (such as LIME), or a convolutional network architecture (such as Expected Gradients and Grad-CAM). ProtoShotXAI only requires a network’s feature layer and data from the classes to construct prototypes. Further, comparable feature attribution maps are created using a simple pixel perturbation which could be enhanced with more intelligent perturbation methods.

Although not explicitly presented in this work, we also attempted different variants of ProtoShotXAI. First, we found that the classification weights are necessary for negative feature attributions because the output of the feature layer often consists of ReLU activations in state-of-the-art neural networks. ReLU activations are either positive or zero, meaning the cosine distance at the feature layer without the classification weights will not produce strong negative feature attributions. In contrast, the Euclidean distance can be used as the similarity score; however, the normalization of cosine distance provided better interpretability across features and was more consistent with the original model response in experiments such as the revolving six. Additionally, the distance between the prototype and query sample can be computed at the output of the final classification layer rather than the weighted feature layer. We found this also produced less favorable results. More specifically, it is common for the final classification layer to have lower dimensionality than the feature layer and a softmax activation function, rather than an unbounded ReLU activation at the feature layer. Both the lower dimensionality and bounded activation functions limited the feature attribution interpretability at the classification layer. The class weighted feature layer produced the best results in our experiments. We also investigated training models to enhance our approach. For example, we implemented an augmentation to the baseline classifier with a few-shot training paradigm. The training objective was to learn adaptive weights between the prototype features and query features such that it would minimize the distance between like classes and maximize the distance between unlike classes. Although the model successfully trained the dynamic weights, it only produced marginally better results. We did not include it here because the benefit of the augmented model was outweighed by the complexity and time cost of having an extra training process. Conveniently, ProtoShotXAI, as presented in this work, does not require any additional training.

6. Conclusions

Machine learning, particularly deep learning, has been shown to provide state-of-the-art performances for numerous benchmarks; however, the explainability of these models and how they reach a decision is a challenging task that limits their utility in some applications. Consequently, it is nearly impossible to quantitatively assess ML approaches’ performance to explain black-box models directly. Instead, indirect qualitative and quantitative measures are used for XAI algorithm performance assessment. Qualitative measures can produce misleading results by exploiting an observer’s confirmation bias and indirect quantitative measures are imperfect assessments. As a result, there is no one-size-fits-all XAI approach and the best approach is dependent on the data and application. It was against this background that we proposed the ProtoShotXAI architecture. ProtoShotXAI is an approach to aid ML practitioners in the exploration of their black-box model. We demonstrated on three datasets that ProtoShotXAI can create feature attribution maps that are quali-

tatively comparable, and arguably superior, to state-of-the-art methods. This work also showed how ProtoShotXAI can be extended to higher-level feature analysis by examining the similarity scores as a function of complex changes to the input as demonstrated in the MNIST revolving **six** and Omniglot sub-character removal experiments. We demonstrate both quantitatively and qualitatively that ProtoShotXAI provides more flexibility for feature analysis and model exploration. Further, the components of the similarity score provide additional insight into model behavior, as demonstrated in the adversarial MNIST experiment. ProtoShotXAI in these specific experiments demonstrates that a user can apply ProtoShotXAI as an exploration tool which will aid in improving the trust in black-box models.

The data used in this work was instrumental in understanding how XAI can be used and evaluated in practice. The digit **eight** was selected due to the convenient presence of other digits within the trace of an **eight**. The image containing the tank and french bulldog was created as a test image with two and only two classes, and all of the other data manipulations (i.e., revolving **six**, sub-character analysis in Omniglot, adversarial analysis) were done intelligently. The intelligent manipulations allowed us to make apriori assumptions about a model’s outcome that could be validated against the information provided by ProtoShotXAI. Our future work aims to expand the intelligent data manipulations done in this work into more comprehensive datasets to better evaluate XAI approaches. Additionally, ProtoShotXAI is applicable to classification networks for many different data types. Although the scope of this work was limited to images, future work will extend ProtoShotXAI to classification networks with sound, text, and data signals.

Acknowledgements

This work was supported by grants from the Department of Energy #DE-NA0003946, Army Research Lab W56KGU-20-C-0002, the National Science Foundation (NSF) CAREER #1943552, and NSF #1950359. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

References

- Amina Adadi and Mohammed Berrada. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160, 2018.
- Michael A Alcorn, Qi Li, Zhitao Gong, Chengfei Wang, Long Mai, Wei-Shinn Ku, and Anh Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4845–4854, 2019.
- Antreas Antoniou, Amos Storkey, and Harrison Edwards. Data augmentation generative adversarial networks. *International Conference on Learning Representations*, 2018.
- Jacob Bien and Robert Tibshirani. Prototype selection for interpretable classification. *The Annals of Applied Statistics*, 5(4):2403–2424, 2011.

- Francesco Bodria, Fosca Giannotti, Riccardo Guidotti, Francesca Naretto, Dino Pedreschi, and Salvatore Rinzivillo. Benchmarking and survey of explanation methods for black box models. *arXiv preprint arXiv:2102.13076*, 2021.
- Steven Bramhall, Hayley Horn, Michael Tieu, and Nibhrat Lohia. Qlime-a quadratic local interpretable model-agnostic explanation approach. *SMU Data Science Review*, 3(1):4, 2020.
- Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a” siamese” time delay neural network. In *Advances in Neural Information Processing Systems*, pages 737–744, 1994.
- Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 839–847. IEEE, 2018.
- Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *International Conference on Learning Representations*, 2019.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
- Radwa ElShawi, Youssef Sherif, Mouaz Al-Mallah, and Sherif Sakr. Ilime: Local and global interpretable model-agnostic explainer of black-box decision. In *European Conference on Advances in Databases and Information Systems*, pages 53–68. Springer, 2019.
- Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott M Lundberg, and Su-In Lee. Improving performance of deep learning models with axiomatic attribution priors and expected gradients. *Nature Machine Intelligence*, pages 1–12, 2021.
- Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- Michael Fink. Object classification from a single example utilizing class relevance metrics. *Advances in neural information processing systems*, 17:449–456, 2004.
- Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In *Neural Information Processing Systems (NeurIPS)*, 2018.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2014.

- Karthik S Gurumoorthy, Amit Dhurandhar, Guillermo Cecchi, and Charu Aggarwal. Efficient data representation by selecting prototypes with importance weights. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 260–269. IEEE, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE CVPR*, pages 770–778, 2016.
- Jeya Vikranth Jeyakumar, Joseph Noor, Yu-Hsi Cheng, Luis Garcia, and Mani Srivastava. How can i explain this to you? an empirical study of deep neural network explanation methods. *Advances in Neural Information Processing Systems*, 2020.
- Been Kim, Rajiv Khanna, and Oluwasanmi O Koyejo. Examples are not enough, learn to criticize! criticism for interpretability. *Advances in neural information processing systems*, 29, 2016.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.
- Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML Deep Learn. Wrksp*, 2015.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, pages 1885–1894. PMLR, 2017.
- Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum. One shot learning of simple visual concepts. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 33, 2011.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*, pages 4768–4777, 2017.
- Erwan Le Merrer and Gilles Trédan. Remote explainability faces the bouncer problem. *Nature Machine Intelligence*, 2:529–539, 2020.
- Nicolas Papernot and Patrick McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *arXiv preprint arXiv:1803.04765*, 2018.
- Tomi Peltola. Local interpretable model-agnostic explanations of bayesian predictive models via kullback-leibler projections. *2nd Workshop on Explainable Artificial Intelligence*, 2018.
- Han Qiu, Yi Zeng, Qinkai Zheng, Shangwei Guo, Tianwei Zhang, and Hewu Li. An Efficient Preprocessing-based Approach to Mitigate Advanced Adversarial Attacks. *IEEE Transactions on Computers*, 2021.
- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. *International Conference on Learning Representations*, 2017.

- 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, 2016.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015a. doi: 10.1007/s11263-015-0816-y.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3): 211–252, 2015b.
- Koosha Sadeghi, Ayan Banerjee, and Sandeep Gupta. A System-Driven Taxonomy of Attacks and Defenses in Adversarial Machine Learning. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 4(4):450–467, 2020.
- Tyler R Scott, Karl Ridgeway, and Michael C Mozer. Adapted deep embeddings: A synthesis of methods for k -shot inductive transfer learning. *arXiv preprint arXiv:1805.08402*, 2018.
- 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. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- Sharath M Shankaranarayana and Davor Runje. Alime: Autoencoder based approach for local interpretability. In *International conference on intelligent data engineering and automated learning*, pages 454–463. Springer, 2019.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and others. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 4080–4090, 2017.

- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328. PMLR, 2017.
- Richard Tomsett, Dan Harborne, Supriyo Chakraborty, Prudhvi Gurram, and Alun Preece. Sanity checks for saliency metrics. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):6021–6029, Apr. 2020.
- Aaron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv:1609.03499*, 2016.
- Oriol Vinyals, Charles Blundell, Tim Lillicrap, Daan Wierstra, and others. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, pages 3630–3638, 2016.
- Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan. Low-shot learning from imaginary data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7278–7286, 2018.
- Meng Ye and Yuhong Guo. Deep triplet ranking networks for one-shot recognition. *arXiv preprint arXiv:1804.07275*, 2018.
- Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian model-agnostic meta-learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 7343–7353, 2018.
- Muhammad Rehman Zafar and Naimul Mefraz Khan. Dlime: A deterministic local interpretable model-agnostic explanations approach for computer-aided diagnosis systems. In *In proceeding of ACM SIGKDD Workshop on Explainable AI/ML (XAI) for Accountability, Fairness, and Transparency*. ACM, 2019.