

# Deep Learning for Case-based Reasoning through Prototypes: A Neural Network that Explains Its Predictions

Oscar Li\*  
Duke University  
runliang.li@duke.edu

Hao Liu\*  
Nanjing University  
141242059@smail.nju.edu.cn

Chaofan Chen  
Duke University  
cfchen@cs.duke.edu

Cynthia Rudin  
Duke University  
cynthia@cs.duke.edu

## Abstract

Deep neural networks are widely used for classification. These deep models often suffer from a lack of interpretability – they are particularly difficult to understand because of their non-linear nature. As a result, neural networks are often treated as “black box” models, and in the past, have been trained purely to optimize the accuracy of predictions. In this work, we create a novel network architecture for deep learning that naturally explains its own reasoning for each prediction. This architecture contains an autoencoder and a special *prototype layer*, where each unit of that layer stores a weight vector that resembles an encoded training input. The encoder of the autoencoder allows us to do comparisons within the latent space, while the decoder allows us to visualize the learned prototypes. The training objective has four terms: an accuracy term, a term that encourages every prototype to be similar to at least one encoded input, a term that encourages every encoded input to be close to at least one prototype, and a term that encourages faithful reconstruction by the autoencoder. The distances computed in the prototype layer are used as part of the classification process. Since the prototypes are learned during training, the learned network naturally comes with explanations for each prediction, and the explanations are loyal to what the network actually computes.

## 1 Introduction

As machine learning algorithms have gained importance for important societal questions, interpretability (transparency) has become a key issue for whether we can trust predictions coming from these models. There have been cases where incorrect data fed into black box models have gone unnoticed, leading to unfairly long prison sentences (e.g., prisoner Glen Rodriguez was denied parole due to an incorrect COMPAS score, Wexler, 2017). In radiology, lack of transparency causes challenges to FDA approval for deep learning products. Because of these issues, “opening the black box” of neural networks has become a debated issue in the media (Citron 2016; Smith 2016; Angwin et al. 2016; Westervelt 2017). Artificial neural networks are particularly

difficult to understand because their highly nonlinear functions do not naturally lend to an explanation that humans are able to process.

In this work, we create an architecture for deep learning that explains its own reasoning process. The learned models naturally come with explanations for each prediction, and the explanations are loyal to what the network actually computes. As we discuss shortly, creating the architecture to encode its own explanations is in contrast with creating explanations for previously trained black box models, and aligns more closely with work on prototype classification and case-based reasoning.

In the past, neural networks have often been designed purely for accuracy, with *posthoc* interpretability analysis. In this case, the network architecture was chosen first, and afterwards one aims to interpret the trained model or the learned high level features. To do the interpretability analysis requires a separate modeling effort. One problem with generating explanations posthoc is that the explanations themselves can change based on the model for the explanation. For instance, it may be easy to create multiple conflicting yet convincing explanations for how the network would classify a single object, none of which are the correct reason for why the object was classified that way. A related issue is that posthoc methods often create explanations that do not make sense to humans. This means that extra modeling is needed to ensure that the explanations are interpretable. This happens, for instance, in the Activation Maximization (AM) approach, where one aims to find an input pattern that produces a maximum model response for a quantity of interest to the user (Erhan et al. 2009). Since the images from AM are not generally interpretable (they tend to be gray), regularized optimization is used to find an interpretable high activation image (Hinton 2012; Lee et al. 2009; van den Oord, Kalchbrenner, and Kavukcuoglu 2016; Nguyen et al. 2016a). When we add regularization, however, the result is a combination of what the network actually computes and the extrinsic regularization. Given that the explanations themselves come from a separate modeling process with strong priors that are not part of training, we then wonder how we can trust the explanations from the posthoc analysis. In fact there is a growing literature discussing the issues mentioned above for AM (Montavon, Samek, and Müller 2017).

\*Contributed equally

For images, posthoc analysis often involves visualization of layers of a neural network. For instance, an alternative to AM is that of Zeiler and Fergus, 2014, who use deconvolution as a technique to visualize what a convolutional neural network (CNN) has learned. Deconvolution is one method for decoding; our method can use any type of decoder to visualize the prototypes, including deconvolution. In addition, we require the other pieces of our architecture in order to create models that are interpretable. We do not consider posthoc analysis in this work.

There is no natural adaptation of the posthoc methods to training a network. Those methods optimize a quantity on the pixel space, and since the resulting models are not generally interpretable, they often include a regularization term such as the density of observations  $p(\mathbf{x})$  to gain interpretability. For training a network, since  $p(\mathbf{x})$  is not a function of the network parameters, adding it into our training objective would not change training.

Our network is a form of *prototype classifier*, where observations are classified based on their proximity to a prototype observation within the dataset. For instance, in our handwritten digit example, we can determine that an observation was classified as a “3” because the network thinks it looks like a particular prototypical “3” within the training set. If the prediction is uncertain, it would identify prototypes similar to the observation from different classes, e.g., “4” is often hard to distinguish from “9”, so we would expect to see prototypes of classes 4 and 9 identified when the network is asked to classify an image of a 9.

Our work is closely aligned with other prototype classification techniques in machine learning (Bien and Tibshirani 2011; Kim, Rudin, and Shah 2014; Priebe et al. 2003; Wu and Tabak 2017). Prototype classification is a classical form of case-based reasoning (Kolodner 1992); however, because our work uses neural networks, the distance measure between prototypes and observations is measured in a flexible latent space. The fact that the latent space is adaptive is the driving force behind its high quality performance.

The word “prototype” is overloaded and has various meanings. For us, a prototype is very close or identical to an observation in the training set, and the set of prototypes are representative of the whole data set. In other contexts, a prototype is not required to be close to any one of the training examples, and could be just a convex combination of several observations. For instance, in the zero-shot and few-shot learning field, prototypes are points in the feature space used to represent a single class, and distance to the prototype determines how an observation is classified. For example, ProtoNets (Snell, Swersky, and Zemel 2017) utilize the mean of several embedded “support” examples as the prototype for each class in few-shot learning. Li and Wang, 2017, use a generative probabilistic model to generate prototypes for zero shot learning, which are points in the feature space. In both cases, prototypes are not optimized to resemble actual observations, and are not required to be interpretable (meaning that their visualizations will not generally resemble natural images), and each class can have only one prototype.

Our deep architecture induces a latent low-dimensional space, and distances to prototypes are computed in that la-

tent space. Using a latent space for distance computation enables us to find a better dissimilarity measure than  $L^2$  on the pixel space. Other works also use latent spaces, e.g., Salakhutdinov and Hinton (2007) use the latent space of a restricted Boltzman machine autoencoder for improved classification, although not for the aim of interpretability.

We use an autoencoder to create our latent low-dimensional space. Autoencoders can learn useful representations and are widely used in representation learning. Autoencoders can be used for clustering (Xie, Girshick, and Farhadi 2015), generative modeling (Kingma and Welling 2014), discovering cross domain relations (Kim et al. 2017), and other tasks.

## 2 Methodology

### 2.1 Network Architecture

Let  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  be the training dataset with  $\mathbf{x}_i \in \mathbb{R}^p$  and  $y_i \in \{1, \dots, K\}$  for each  $i \in \{1, \dots, n\}$ . Our model architecture consists of two components: an autoencoder (including an encoder,  $f: \mathbb{R}^p \rightarrow \mathbb{R}^q$ , and a decoder,  $g: \mathbb{R}^q \rightarrow \mathbb{R}^p$ ) and a prototype classification network  $h: \mathbb{R}^q \rightarrow \mathbb{R}^K$ , illustrated in Figure 1. The network uses the autoencoder to reduce the dimensionality of the input and learn useful features for prediction; then it uses the encoded input to produce a probability distribution over the  $K$  classes through the prototype classification network  $h$ . The network  $h$  is made up of three layers: a prototype layer,  $p: \mathbb{R}^q \rightarrow \mathbb{R}^m$ , a fully-connected layer  $w: \mathbb{R}^m \rightarrow \mathbb{R}^K$ , and a softmax layer,  $s: \mathbb{R}^K \rightarrow \mathbb{R}^K$ . The network learns  $m$  prototype vectors  $\mathbf{p}_1, \dots, \mathbf{p}_m \in \mathbb{R}^q$  (each corresponds to a *prototype unit* in the architecture) in the latent space. The prototype layer  $p$  computes the squared  $L^2$  distance between the encoded input  $\mathbf{z} = f(\mathbf{x}_i)$  and each of the prototype vectors:

$$p(\mathbf{z}) = [\|\mathbf{z} - \mathbf{p}_1\|_2^2, \|\mathbf{z} - \mathbf{p}_2\|_2^2, \dots, \|\mathbf{z} - \mathbf{p}_m\|_2^2]^\top. \quad (1)$$

In Figure 1, the *prototype unit* corresponding to  $\mathbf{p}_j$  executes the computation  $\|\mathbf{z} - \mathbf{p}_j\|_2^2$ . The fully-connected layer  $w$  computes weighted sums of these distances  $Wp(\mathbf{z})$ , where  $W$  is a  $K \times m$  weight matrix. These weighted sums are then normalized by the softmax layer  $s$  to output a probability distribution over the  $K$  classes. The  $k$ -th component of the output of the softmax layer  $s$  is defined by

$$s(\mathbf{v})_k = \frac{\exp(v_k)}{\sum_{k'=1}^K \exp(v_{k'})} \quad (2)$$

where  $v_k$  is the  $k$ -th component of the vector  $\mathbf{v} = Wp(\mathbf{z}) \in \mathbb{R}^K$ .

During prediction, the model outputs the class that it thinks is the most possible. In essence, our classification algorithm is distance-based on the low-dimensional learned feature space. A special case is when we use one prototype for every class (let  $m = K$ ) and set the weight matrix of the fully-connected layer to the negative identity matrix,  $W = -I_{K \times K}$  (i.e.  $W$  is not learned during training). Then the data will be predicted to be in the same class as the nearest prototype in the latent space. More realistically, we typically do not know how many prototypes should be assigned

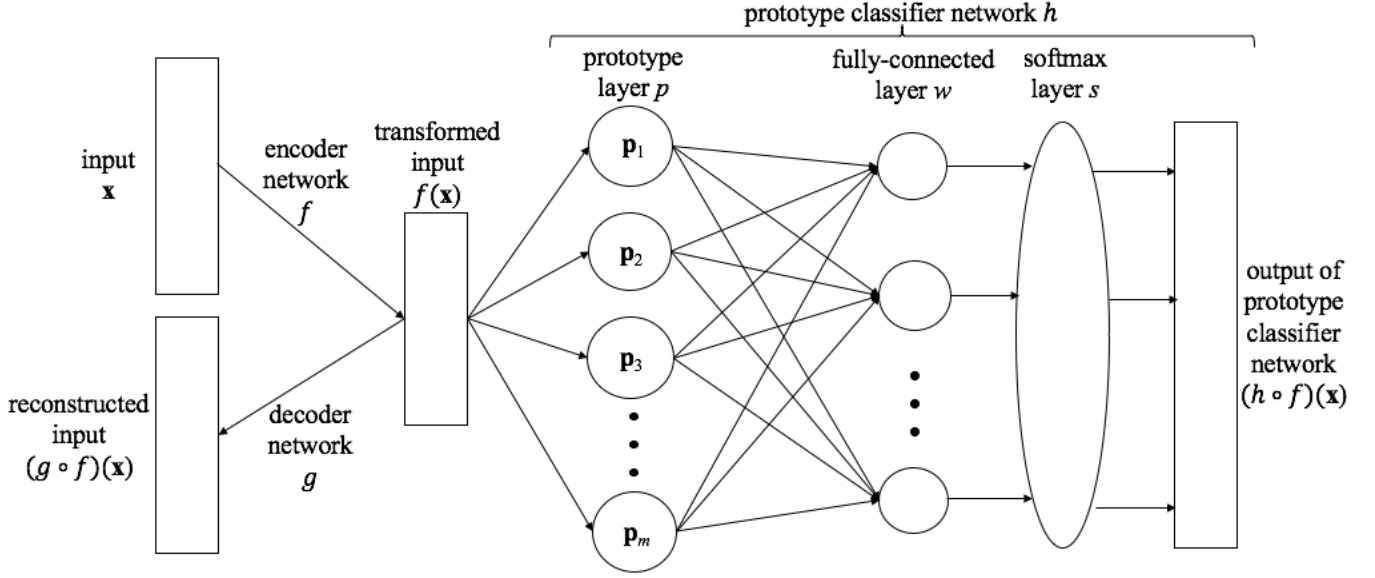


Figure 1: Network Architecture

to each class, and we may want a different number of prototypes from the number of classes, i.e.,  $m \neq K$ . In this case, we allow  $W$  to be learned by the network, and, as a result, the distances to all the prototype vectors will contribute to the probability prediction for each class.

This network architecture has at least three advantages. First, unlike traditional case-based learning methods, the new method automatically learns useful features. For image datasets, which have dimensions equal to the number of pixels, if we perform classification using the original input space or use hand-crafted feature spaces, the methods tend to perform poorly (e.g.,  $k$ -Nearest Neighbors). Second, because the prototype vectors live in the same space as the encoded inputs, we can feed these vectors into the decoder and visualize the learned prototypes throughout the training process. This property, coupled with the case-based reasoning nature of the prototype classification network  $h$ , gives users the ability to interpret how the network reaches its predictions and visualize the prototype learning throughout the full training process without *posthoc* analysis. Third, when we allow the weight matrix  $W$  to be learnable, we are able to tell from the strengths of the learned weight connections which prototypes are more representative of which class.

## 2.2 Cost Function

The network's cost function reflects the needs of both accuracy and interpretability. In addition to the classification error, there is a (standard) term that penalizes the reconstruction error of the autoencoder. There are two new error terms that encourage the learned prototype vectors to correspond to meaningful points in the input space; in our case studies, these points are realistic images. The full four terms are provided mathematically below.

We use the standard cross-entropy loss for penalizing the misclassification. The cross-entropy loss on the training data

$D$  is denoted by  $E$ , and is given by

$$E(h \circ f, D) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K -\mathbb{1}[y_i = k] \log((h \circ f)_k(\mathbf{x}_i)) \quad (3)$$

where  $(h \circ f)_k$  is the  $k$ -th component of  $(h \circ f)$ . We use the squared  $L^2$  distance between the original and reconstructed input for penalizing the autoencoder's reconstruction error. The reconstruction loss, denoted by  $R$ , on the training data  $D$  is given by

$$R(g \circ f, D) = \frac{1}{n} \sum_{i=1}^n \|(g \circ f)(\mathbf{x}_i) - \mathbf{x}_i\|_2^2. \quad (4)$$

The two interpretability regularization terms are formulated as follows:

$$R_1(\mathbf{p}_1, \dots, \mathbf{p}_m, D) = \frac{1}{m} \sum_{j=1}^m \min_{i \in [1, n]} \|\mathbf{p}_j - f(\mathbf{x}_i)\|_2^2, \quad (5)$$

$$R_2(\mathbf{p}_1, \dots, \mathbf{p}_m, D) = \frac{1}{n} \sum_{i=1}^n \min_{j \in [1, m]} \|f(\mathbf{x}_i) - \mathbf{p}_j\|_2^2. \quad (6)$$

Here both terms are averages of minimum squared distances. The minimization of  $R_1$  would require each prototype vector to be as close as possible to one of the training examples in the latent space. As long as we choose the decoder network to be a continuous function, we should expect two very close vectors in the latent space to be decoded to similar-looking images. Thus,  $R_1$  will push the prototype vectors to have meaningful decodings in the pixel space. The minimization of  $R_2$  would require every encoded training example to be as close as possible to one of the prototype vectors. This means that  $R_2$  will cluster the training examples around prototypes in the latent space. We notice here that although  $R_1$

and  $R_2$  involves a minimization function that is not differentiable everywhere, these terms are differentiable almost everywhere and many modern deep learning libraries support this type of differentiation. Ideally,  $R_1$  would take the minimum distance over the entire training set for every prototype; therefore, the gradient computation would grow linearly with the size of the training set. However, this would be impractical during optimization for a large dataset. To address this problem, we relax the minimization to be over only the random minibatch used by the Stochastic Gradient Descent (SGD) algorithm. For the other three terms, since each of them is a summation over the entire training set, it is natural to apply SGD to randomly selected batches for gradient computation.

Putting everything together, the cost function, denoted by  $L$ , on the training data  $D$  with which we train our network  $(f, g, h)$ , is given by

$$\begin{aligned} L((f, g, h), D) = & E(h \circ f, D) + \lambda R(g \circ f, D) \\ & + \lambda_1 R_1(\mathbf{p}_1, \dots, \mathbf{p}_m, D) \\ & + \lambda_2 R_2(\mathbf{p}_1, \dots, \mathbf{p}_m, D), \end{aligned} \quad (7)$$

where  $\lambda$ ,  $\lambda_1$ , and  $\lambda_2$  are real-valued hyperparameters that adjust the ratios between the terms.

### 3 Case Study 1: Handwritten Digits

We now begin a detailed walkthrough of applying our model to the well-known MNIST dataset. The Modified NIST Set (MNIST) is a benchmark dataset of gray-scale images of segmented and centered handwritten digits (Lecun et al., 1998). We used 55,000 training examples, 5,000 validation examples, and 10,000 testing examples, where every image is of size  $28 \times 28$  pixels. We preprocess the images so that every pixel value is in  $[0, 1]$ . This section is organized as follows: we first introduce the architecture and the training details, then compare the performance of our network model with other noninterpretable network models (including a regular convolutional neural network), and finally visualize the learned prototypes, the weight matrix  $W$ , and how a specific image is classified.

#### 3.1 Architecture Details

Hinton and Salakhutdinov (2006) showed that a multilayer fully connected autoencoder network can achieve good reconstruction on MNIST even when using a very low dimensional latent space. We chose a multilayer convolutional autoencoder with a symmetric architecture for the encoder and decoder to be our model’s autoencoder; these types of networks tend to reduce spatial feature extraction redundancy on image data sets and learn useful hierarchical features for producing state-of-the-art classification results. Each convolutional layer consists of a convolution operation followed by a pointwise nonlinearity. We achieve down-sampling in the encoder through strided convolution, and use strided deconvolution in the corresponding layer of the decoder. After passing the original image through the encoder, the network flattens the resulted feature maps into a code vector and feeds it into the prototype layer. The resulting unflattened feature maps are fed into the decoder to reconstruct the

original image. To visualize a prototype vector in the pixel space, we first reshape the vector to be in the same shape as the encoder output and then feed the shaped vector (now a series of feature maps) into the decoder.

The autoencoder in our network has four convolutional layers in both the encoder and decoder. All four convolutional layers in the encoder use kernels of size  $3 \times 3$ , same zero padding, and stride of size 2 in the convolution stage. The filters in the corresponding layers in the encoder and decoder are not constrained to be transposes of each other. Each of the outputs of the first three layers has 32 feature maps, while the last layer has 10. Given an input image of dimension  $28 \times 28 \times 1$ , the shape of the encoder layers are thus:  $14 \times 14 \times 32$ ;  $7 \times 7 \times 32$ ;  $4 \times 4 \times 32$ ;  $2 \times 2 \times 10$ , and therefore the network compresses every 784-dimensional image input to a 40-dimensional code vector ( $2 \times 2 \times 10$ ). Every layer uses the sigmoid function  $f(x) = \frac{1}{1+e^{-x}}$  as the nonlinear transformation. We specifically use the sigmoid function in the last encoder layer so that the output of the encoder is restricted to the unit hypercube  $(0, 1)^{40}$ . This allows us to initialize 15 prototype vectors uniformly at random in that hypercube. We do not use the rectified linear unit (ReLU – Krizhevsky, Sutskever, and Hinton, 2012) because if ReLU were used in the last encoder layer, it would be more difficult to initialize the prototype vectors, as initial states throughout  $\mathbb{R}_{\geq 0}^{40}$  would need to be explored, and the network would take longer to stabilize. We also specifically choose the sigmoid function for the last decoder layer to make the range of pixel values in the reconstructed output  $(0, 1)$ , roughly the same as the preprocessed image’s pixel range.

#### 3.2 Training Details

We set all the hyperparameters  $\lambda$ ,  $\lambda_1$ ,  $\lambda_2$  to 0.05 and the learning rate to 0.0001. We minimized (7) as a whole: we did not employ a greedy stepwise optimization for different layers of the autoencoder nor did we first train the autoencoder and then the prototype classification network.

Our goal in this work is not just to obtain reasonable accuracy, but also interpretability. We used only a few of the possible techniques for improving performance in neural networks generally, and it is possible that using more techniques would improve accuracy (see supplement for more details). In particular, we used the data augmentation technique *elastic deformation* (Simard, Steinkraus, and Platt 2003) to improve the prediction accuracy and reduce potential overfitting. The set of all elastic deformations is a superset of affine transformations. For every mini-batch of size 250 that we randomly sampled from the training set, we applied a random elastic distortion where a Gaussian filter of standard deviation equal to 4 and a scaling factor of 20 were used for the displacement field. Due to the randomness in the data augmentation process, the network sees a slightly different set of images during every epoch, which significantly reduces overfitting.

#### 3.3 Accuracy

After training for 1500 epochs, our model achieved a classification accuracy of 99.53% on the standard MNIST training set and 99.22% on the standard MNIST test set.

To examine how the two key elements of the interpretable network (the autoencoder and prototype layer) affect predictive power, we performed a type of ablation study. In particular, we trained two classification networks that are similar to ours, but removed some key pieces in each of the networks. The first network substitutes the prototype layer with a fully-connected layer whose output is a 15-dimensional vector, the same dimension as the output from the prototype layer; the second network also removes the decoder and changes the nonlinearity to ReLU. The second network is just a regular convolutional neural network that has similar architectural complexity to LeNet 5 introduced in Lecun et al. (1998). After training both networks using elastic deformation for 1500 epochs, we obtained test accuracies of 99.24% and 99.23% respectively. These test accuracies, along with the test accuracy of 99.2% reported in Lecun et al. (1998), are comparable to the test accuracy of 99.22% obtained using our interpretable network. This result demonstrates that changing from a traditional convolutional neural network to our interpretable network architecture does not hinder the predictive ability of the network (at least not in this case).

In general, it is not always true that accuracy needs to be sacrificed to obtain interpretability; there could be many models that are almost equally accurate. The extra terms in the cost function (and changes in architecture) encourage the model to be more interpretable among the set of approximately equally accurate models.

### 3.4 Visualization

Let us first discuss the quality of the autoencoder, because good performance of the autoencoder will allow us to interpret the prototypes. After training, our network’s autoencoder achieved an average squared  $L^2$  reconstruction error of 4.22 over the undeformed training set, where examples are shown in Figure 2. This reconstruction result assures us that the decoder can faithfully map the prototype vectors to the pixel space.

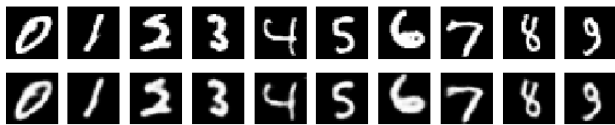


Figure 2: Some random images from the training set in the first row and their corresponding reconstructions in the second row.



Figure 3: The 15 prototype vectors of MNIST visualized in the pixel space.

We sent the learned prototype vectors through the decoder and visualized them in Figure 3. The decoded prototype images are sharp-looking and mostly resemble real-life hand-written digits, owing to the interpretability terms  $R_1$  and  $R_2$  in the cost function. Note that there is not a one-to-one correspondence between classes and prototypes. Since we multiply the output of the prototype layer by a learnable weight matrix prior to feeding it into the softmax layer, the distances from an encoded image to each prototype have differing effects on the predicted class probabilities.

We now look at the transposed weight matrix connecting the prototype layer to the softmax layer, shown in Table 1, to see the influence of the distance to each prototype on every class. We observe that each decoded prototype is visually similar to an image of the class for which the corresponding entry in the weight matrix has a significantly negative value. We will call the class to which a decoded prototype is visually similar the *visual class* of the prototype.

The reason for such a significantly negative value can be understood as follows. The prototype layer is computing the dissimilarity between an input image and a prototype through the squared  $L_2$  distance between their representations in the latent space. This means that the farther away an encoded image  $f(\mathbf{x}_i)$  is to a specific prototype  $\mathbf{p}_j$ , the less the network will think that  $\mathbf{x}_i$  is in the visual class of the prototype ( $\mathbf{p}_j$ ). Thus as soon as  $\|\mathbf{p}_j - f(\mathbf{x}_i)\|_2^2$  becomes large, the probability of  $\mathbf{x}_i$  being in the visual class of  $\mathbf{p}_j$  should dramatically decrease. Therefore, the weight connection between a prototype and its visual class is highly negative. Conversely, if  $\mathbf{x}_i$  is in the visual class of the prototype  $\mathbf{p}_j$ , even if the weight is highly negative, it is multiplied by a very small distance.

An interesting prototype learned by the network is the last prototype in Table 1. It is visually similar to an image of class 2; however, it has strong negative weight connections with class 7 and class 8 as well. Therefore, we can think of this prototype as being shared by these three classes, which means that an encoded input image that is far away from this prototype in latent space would be unlikely to be an image of 7, 8, or 2. This should not be too surprising: if we look at this decoded prototype image carefully, we can see that if we hide the tail of the digit, it would look like an image of 7; if we connect the upper-left endpoint with the lower-right endpoint, it would look like an image of 8.

Let us now look at the learned prototypes in Figure 3. The three prototypes for class 6 seem to represent different writing habits in terms of what the loop of 6 looks like. The first and third 6’s have their loops end at the bottom while the second 6’s loop ends more on the side. The 2’s show similar variation. As for the two 3’s, the two prototypes corresponds to different curvatures.

Let us look into the model as it produces a prediction for a specific image of digit 6, shown on the left of Table 2. The distances computed by the prototype layer between the encoded input image and each of the prototypes are shown below the decoded prototypes in Table 2, and the three smallest distances correspond to the three prototypes that resemble 6 after decoding. We observe here that these three distances are different from each other, and the encoded input image is

	0	1	2	3	4	5	6	7	8	9
8	-0.07	7.77	1.81	0.66	4.01	2.08	3.11	4.10	-20.45	-2.34
9	2.84	3.29	1.16	1.80	-1.05	4.36	4.40	-0.71	0.97	-18.10
0	-25.66	4.32	-0.23	6.16	1.60	0.94	1.82	1.56	3.98	-1.77
7	-1.22	1.64	3.64	4.04	0.82	0.16	2.44	-22.36	4.04	1.78
3	2.72	-0.27	-0.49	-12.00	2.25	-3.14	2.49	3.96	5.72	-1.62
6	-5.52	1.42	2.36	1.48	0.16	0.43	-11.12	2.41	1.43	1.25
3	4.77	2.02	2.21	-13.64	3.52	-1.32	3.01	0.18	-0.56	-1.49
1	0.52	-24.16	2.15	2.63	-0.09	2.25	0.71	0.59	3.06	2.00
6	0.56	-1.28	1.83	-0.53	-0.98	-0.97	-10.56	4.27	1.35	4.04
6	-0.18	1.68	0.88	2.60	-0.11	-3.29	-11.20	2.76	0.52	0.75
5	5.98	0.64	4.77	-1.43	3.13	-17.53	1.17	1.08	-2.27	0.78
2	1.53	-5.63	-8.78	0.10	1.56	3.08	0.43	-0.36	1.69	3.49
2	1.71	1.49	-13.31	-0.69	-0.38	4.55	1.72	1.59	3.18	2.19
4	5.06	-0.03	0.96	4.35	-21.75	4.25	1.42	-1.27	1.64	0.78
2	-1.31	-0.62	-2.69	0.96	2.36	2.83	2.76	-4.82	-4.14	4.95

Table 1: Transposed weight matrix (every entry rounded off to 2 decimal places) between the prototype layer and the softmax layer. Each row represents a prototype node whose decoded image is shown in the first column. Each column represents a digit class. The most negative weight is shaded for each prototype. In general, for each prototype, its most negative weight is towards its visual class except for the prototype in the last row.

significantly closer to the third “6” prototype than the other two. This indicates that our model is indeed capturing the subtle differences within the same class.

After the prototype layer computes the 15-dimensional vector of distances shown in Table 2, it is multiplied by the weight matrix in Table 1, and the output is the unnormalized probability vector used as the logit for the softmax layer. The predicted probability of class 6 for this specific image is 99.99%.

	8	9	0	7	3
6	0.98	1.47	0.70	1.55	1.49
6	0.29	1.69	1.02	0.41	0.15
5	0.88	1.40	1.45	1.28	1.28

Table 2: The rounded-off distances between a test image 6 and every prototype in the latent space.

## 4 Case Study 2: Cars

The second dataset we use consists of rendered images, each with  $64 \times 64 \times 3$  pixels, of 3D car models with varying azimuth angles at  $15^\circ$  intervals, from  $-75^\circ$  to  $75^\circ$  (Fidler, Dickinson, and Urtasun 2012). There are 11 views of each car and every car’s class label is one of the 11 angles. The dataset is split into a training set ( $169 \times 11 = 1859$  images) and a test set ( $14 \times 11 = 154$  images).

We use two convolutional layers in both the encoder and decoder. The first and the second layer in the encoder uses respectively 32 and 10 convolutional filters of size  $5 \times 5$ ,



Figure 4: Three cars at 11 angles from car dataset.

stride 2, and no zero padding. The architecture of the decoder is symmetric to that of the encoder. We use the sigmoid activation function in the last layer of the encoder and the decoder, and leaky ReLU in all other autoencoder layers. We set the number of our prototypes to be eleven, which is the same as the number of classes. Figure 5 shows the eleven decoded prototypes from our model. The network has determined that the color of a car is not important in determining the angle, so all of the decoded prototypes are of the same “average” color. The learned weight matrix  $W$  is shown in Table 4 in the Supplementary Material. We compared our model with a network without the interpretable parts, in which we removed the decoder and replaced the prototype layer with a fully connected layer of the same size. The accuracies for these two models are shown in Table 3. The result again illustrates that we do not sacrifice much accuracy when including the interpretability elements into the network.

	interpretable	non-interpretable
train acc	98.2%	99.8%
test acc	93.5%	94.2%

Table 3: Car dataset accuracy.



Figure 5: Decoded prototypes when we include  $R_1$  and  $R_2$ .

We use this case study to illustrate the importance of the two interpretability terms  $R_1$  and  $R_2$  in our cost function. If we remove both  $R_1$  and  $R_2$ , the decoded prototypes will not look like real images, as shown in Figure 6. If we leave out only  $R_1$ , the decoded prototypes will again not look like real observations, as shown in Figure 7. If we remove only  $R_2$ , the network chooses prototypes that do not fully represent the input space, and some of the prototypes tend to be similar to each other, as shown in Figure 8. Intuitively,  $R_1$  pushes every prototype to be close to a training example in the latent space so that the decoded prototypes can be realistic, while  $R_2$  forces every training example to find a close prototype in the latent space, thereby encouraging the prototypes to spread out over the entire latent space and to be distinct from each other. In other words,  $R_1$  helps make the prototypes meaningful, and  $R_2$  keeps the explanations faithful in forcing the network to use nearby prototypes for classification.



Figure 6: Decoded prototypes when we remove  $R_1$  and  $R_2$ .

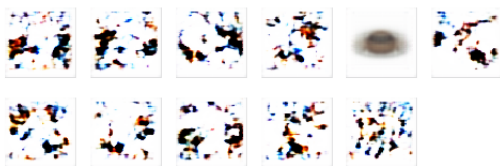


Figure 7: Decoded prototypes when we remove  $R_1$ .

### 5 Case Study 3: Fashion MNIST

Fashion MNIST (Xiao, Rasul, and Vollgraf 2017) is a dataset of Zalando’s article images, consisting of a training set of 60,000 examples and a test set of 10,000 examples. Each example is a  $28 \times 28$  grayscale image, associated with a label from 10 classes, each being a type of clothes item. The dataset shares the same image size and structure of training and testing splits as MNIST.



Figure 8: Decoded prototypes when we remove  $R_2$ .

We ran the same model from Case Study 1 on this fashion dataset and achieved a testing accuracy of 89.95%. This result is comparable to those obtained using standard convolutional neural networks with max pooling reported on the dataset website (87.6-92.5% for networks that use similar architecture complexity as ours, Fashion-MNIST, 2017). The learned prototypes are shown in Figure 9. For each class, there is at least one prototype representing that class. The learned prototypes have fewer details (such as stripes, presence of a collar, texture) than the original images. This again shows that the model has recognized what information is important in this classification task – the contour shape of the input is more useful than its fine-grained details. The learned weight matrix  $W$  is shown in Table 5 in the Supplementary Material.

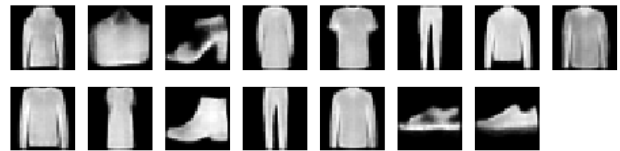


Figure 9: 15 decoded prototypes for Fashion-MNIST.

## 6 Discussion and Conclusion

We combine the strength of deep learning and the interpretability of case-based reasoning to make an interpretable deep neural network model. The prototypes can provide useful insight into the inner workings of the network, the relationship between classes, and the important aspects of the latent space, as demonstrated here. Although our model does not provide a full solution to problems with accountability and transparency of black box decisions, it does allow us to partially trace the path of classification for a new observation.

We noticed in our experiments that the addition of the two interpretability terms  $R_1$  and  $R_2$  tended to act as regularizers, and helped to make the network robust to overfitting. The extent to which interpretability reduces overfitting is a topic that could be explored in future work.

**Code:** Our code is publicly available at this URL: *placeholder*.

**Acknowledgments:** This work sponsored in part by MIT Lincoln Laboratory.



## References

- [2016] Angwin, J.; Larson, J.; Mattu, S.; and Kirchner, L. 2016. Machine bias. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- [2011] Bien, J., and Tibshirani, R. 2011. Prototype selection for interpretable classification. *Annals of Applied Statistics* 5(4):2403–2424.
- [2016] Citron, D. 2016. (Un)fairness of risk scores in criminal sentencing. *Forbes, Tech section*.
- [2009] Erhan, D.; Bengio, Y.; Courville, A.; and Vincent, P. 2009. Visualizing higher-layer features of a deep network. Technical Report 1341, University of Montreal. Also presented at the ICML 2009 Workshop on Learning Feature Hierarchies, Montreal, Canada.
- [2017] Fashion-MNIST. 2017. Github repository website. <https://github.com/zalando-research/fashion-mnist>. Online; accessed September 7, 2017.
- [2012] Fidler, S.; Dickinson, S.; and Urtasun, R. 2012. 3d object detection and viewpoint estimation with a deformable 3d cuboid model. In Pereira, F.; Burges, C. J. C.; Bottou, L.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc. 611–619.
- [2014] Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A. C.; and Bengio, Y. 2014. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 2672–2680.
- [2015] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Deep residual learning for image recognition. *CoRR* abs/1512.03385.
- [2006] Hinton, G. E., and Salakhutdinov, R. R. 2006. Reducing the dimensionality of data with neural networks. *Science* 313(5786):504–507.
- [2012] Hinton, G. E. 2012. A practical guide to training restricted boltzmann machines. In *Neural Networks: Tricks of the Trade - Second Edition*. Springer. 599–619.
- [2016] Huang, G.; Liu, Z.; and Weinberger, K. Q. 2016. Densely connected convolutional networks. *CoRR* abs/1608.06993.
- [2017] Kim, T.; Cha, M.; Kim, H.; Lee, J. K.; and Kim, J. 2017. Learning to discover cross-domain relations with generative adversarial networks. In *Proceedings of International Conference on Machine Learning (ICML)*.
- [2014] Kim, B.; Rudin, C.; and Shah, J. 2014. The bayesian case model: A generative approach for case-based reasoning and prototype classification. In *Proceedings of Neural Information Processing Systems (NIPS)*.
- [2014] Kingma, D. P., and Welling, M. 2014. Auto-encoding variational bayes. In *The International Conference on Learning Representations (ICLR)*.
- [1992] Kolodner, J. 1992. An introduction to case-based reasoning. *AI Review*.
- [2012] Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In Pereira, F.; Burges, C. J. C.; Bottou, L.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc. 1097–1105.
- [1998] Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11):2278–2324.
- [2009] Lee, H.; Grosse, R. B.; Ranganath, R.; and Ng, A. Y. 2009. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, 609–616.
- [2017] Li, Y., and Wang, D. 2017. Zero-shot learning with generative latent prototype model. *CoRR* abs/1705.09474.
- [2017] Montavon, G.; Samek, W.; and Müller, K. 2017. Methods for interpreting and understanding deep neural networks. *CoRR* abs/1706.07979.
- [2016a] Nguyen, A.; Dosovitskiy, A.; Yosinski, J.; Brox, T.; and Clune, J. 2016a. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*, 3387–3395.
- [2016b] Nguyen, A. M.; Dosovitskiy, A.; Yosinski, J.; Brox, T.; and Clune, J. 2016b. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *CoRR* abs/1605.09304.
- [2003] Priebe, C. E.; Marchette, D. J.; DeVinney, J. G.; and Socolinsky, D. A. 2003. Classification using class cover catch digraphs. *Journal of classification* 20(1):003–023.
- [2007] Salakhutdinov, R., and Hinton, G. E. 2007. Learning a nonlinear embedding by preserving class neighbourhood structure. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics, AISTATS*, 412–419.
- [2003] Simard, P. Y.; Steinkraus, D.; and Platt, J. C. 2003. Best practices for convolutional neural networks applied to visual document analysis. In *Proceedings of the Seventh International Conference on Document Analysis and Recognition (ICDAR), Volume 2*.
- [2016] Smith, M. 2016. In wisconsin, a backlash against using data to foretell defendants’ futures. *New York Times*.
- [2017] Snell, J.; Swersky, K.; and Zemel, R. S. 2017. Prototypical networks for few-shot learning. *CoRR* abs/1703.05175.
- [2016] van den Oord, A.; Kalchbrenner, N.; and Kavukcuoglu, K. 2016. Pixel recurrent neural networks. In *Proceedings of the 33rd International Conference on Machine Learning, (ICML)*, 1747–1756.
- [2017] Westervelt, E. 2017. Did a bail reform algorithm contribute to this San Francisco man’s murder? *National Public Radio, Law*.
- [2017] Wexler, R. 2017. When a computer program keeps you in jail: How computers are harming criminal justice. *New York Times*.



- [2017] Wu, C., and Tabak, E. G. 2017. Prototypal analysis and prototypal regression.
- [2017] Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.
- [2015] Xie, J.; Girshick, R. B.; and Farhadi, A. 2015. Un-supervised deep embedding for clustering analysis. *CoRR* abs/1511.06335.
- [2014] Zeiler, M. D., and Fergus, R. 2014. Visualizing and understanding convolutional networks. In *In European Conference on Computer Vision*, 818–833. Springer.

## **A Supplementary Material**

### **A.1 Possible Techniques for Improving the Performance of Our Model**

In our demonstrations we work with a basic network that demonstrates the main idea. If one wants to apply our model to more complex datasets, there are many techniques that can be adopted. We list a few of them below:

- Optimize the structure of the network. For instance, add more layers or using more powerful architectures such as ResNet (He et al. 2015), DenseNet (Huang, Liu, and Weinberger 2016), etc. to improve the accuracy.
- Add stronger regularizers to the training objective to make the prototype more realistic, such as using Generative Adversarial Networks (GAN) (Goodfellow et al. 2014). GAN has been used for improved visualization in posthoc analysis such as AM (Nguyen et al. 2016b).

### **A.2 Weight Matrices**

We show the learned weight matrices of Case Studies 2 and 3 in this section. The most negative weight connections are shaded for each prototype.












	$-75^\circ$	$-60^\circ$	$-45^\circ$	$-30^\circ$	$-15^\circ$	$0^\circ$	$15^\circ$	$30^\circ$	$45^\circ$	$60^\circ$	$75^\circ$
	1.33	0.45	0.01	0.70	-0.33	-2.11	0.12	0.61	-0.35	1.07	1.49
	1.69	-0.48	-2.93	0.07	1.39	0.63	0.86	1.92	0.86	1.09	0.50
	1.03	0.35	0.81	-0.21	0.54	-2.97	0.27	0.67	1.46	0.23	0.42
	1.36	1.15	1.62	0.65	0.39	0.76	0.33	-3.77	-0.13	0.78	1.05
	-0.88	0.46	1.80	1.46	1.55	1.93	1.20	0.80	2.25	-0.83	-3.01
	0.26	0.66	0.93	1.16	1.21	1.41	0.99	1.63	-0.44	-1.33	-1.09
	0.58	1.30	0.21	0.46	-3.93	0.68	0.76	0.68	0.86	0.57	1.42
	-2.09	-1.89	0.52	2.01	1.22	1.21	1.59	1.75	1.06	1.53	0.90
	0.17	1.14	0.83	1.29	1.39	0.99	1.51	0.57	-2.46	-1.26	1.19
	1.10	1.11	0.13	0.59	0.74	0.31	-4.29	0.05	0.05	1.17	0.90
	-0.02	0.26	0.68	-3.42	0.43	1.55	1.51	0.09	1.51	1.66	0.99

Table 4: Transposed weight matrix learned from Case Study 2 with both  $R_1$  and  $R_2$ .




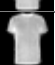
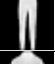




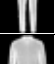





	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot
	-0.72	3.02	-11.00	1.72	-4.62	3.26	9.26	0.32	1.88	0.82
	0.18	4.20	5.24	-1.24	0.08	6.86	1.03	-2.11	-26.31	1.91
	3.16	0.30	2.78	3.69	0.78	-16.90	-3.51	3.12	5.05	-1.28
	4.50	6.88	3.39	-14.51	-2.76	5.59	-1.29	-1.44	-0.62	4.38
	-16.56	6.87	0.32	2.13	6.76	3.02	-5.43	4.58	-0.09	-0.17
	-0.92	-15.36	0.42	1.07	-1.11	-1.21	3.92	0.66	3.63	1.89
	-0.04	0.25	1.97	5.97	-14.13	0.91	0.90	2.36	1.15	3.38
	5.57	0.23	5.69	5.41	-2.03	-0.31	-12.96	2.48	2.76	3.23
	1.27	1.42	-11.92	1.09	7.31	3.07	-3.23	4.47	6.11	-1.23
	2.31	5.92	0.16	-17.12	3.65	5.38	-0.03	-0.53	1.21	-0.04
	-0.32	2.01	-0.97	-1.96	3.47	0.84	3.09	-5.15	-2.68	-24.82
	-0.97	-14.13	0.06	1.68	-1.30	-0.26	0.74	2.57	4.44	2.73
	1.73	1.06	0.40	-0.26	0.14	3.03	-5.41	2.63	-0.86	2.59
	0.15	0.38	2.89	1.21	2.41	-13.93	-1.93	1.10	2.04	0.58
	1.62	0.95	0.12	-2.27	0.28	1.45	4.55	-25.07	-4.31	-1.38

Table 5: Tranposed weight matrix learned in Case Study 3.