
Not Just A Black Box: Interpretable Deep Learning by Propagating Activation Differences

Avanti Shrikumar¹ (avanti@stanford.edu), Peyton Greenside² (pgreens@stanford.edu)
Anna Y. Shcherbina² (annashch@stanford.edu), Anshul Kundaje^{1,3} (akundaje@stanford.edu)

1. Department of Computer Science, Stanford University, CA, USA
2. Biomedical Informatics, Stanford University, CA, USA
3. Department of Genetics, Stanford University, CA, USA

Abstract

A common criticism of neural networks is their lack of interpretability or “black box” nature. This is a major barrier to adoption in application domains such as biology, where interpretability is important. Here we present LIFTPAD (Learning Important Features Through Propagating Activation Differences), an intuitive and effective method for scoring the contributions of inputs to the output of a neural network. LIFTPAD compares the activation of each neuron to its ‘reference activation’, and assigns a contribution score to each of the neuron’s inputs proportionally to how the inputs differ from their own ‘reference activations’. We apply LIFTPAD to models trained on natural images as well as models trained on genomic sequences. We show that LIFTPAD correctly identifies important features in both cases and has significant advantages over gradient-based methods.

1. Introduction

As neural networks become increasingly popular, their reputation as a “black box” presents a barrier to adoption in fields where interpretability is paramount. Understanding the features that lead to a particular output builds trust with users and can lead to novel scientific discoveries. A common approach is to leverage the gradients of a particular output with respect to the individual inputs (Simonyan et al. - cite) - however, such approaches are limited because activation functions such as Rectified Linear Units (ReLU) have a gradient of zero when they are not firing, yet a ReLU that does not fire can still carry information - particularly if

the associated bias is positive (which means that in the absence of any input, the ReLU does fire). Similarly, sigmoid or tanh activations are popular choices for the activation functions of gates in memory units of recurrent neural networks such as GRUs and LSTMs (cite), but these activations have a near-zero gradient at high or low inputs even though such inputs can be very significant.

In this paper, we present a general method for assigning feature importance that focuses on the difference between a neuron’s activation and its ‘reference’ activation, where the reference activation is the activation that the neuron has when the network is provided a ‘reference input’. The reference input is defined on a domain-specific basis according to what is appropriate for the task at hand. This circumvents the aforementioned limitation of gradient-based approaches because, rather than relying on the local gradient at the point of activation, the activation is simply compared to its reference value. In the case of a ReLU that fires when provided the reference input (as is typically the case for ReLUs that have a positive bias), our method assigns a negative ‘difference from reference’ when the ReLU does not fire. Similarly, in the case of a sigmoidal unit which has an output of 0.5 when provided the reference input, our method would assign a positive ‘difference from reference’ when the sigmoidal unit has an output near 1, even though the gradient at that output is negligible.

2. LIFTPAD Method

We denote the contribution of x to y as C_{xy} . Let the activation of a neuron n be denoted as A_n . Further, let the *reference* activation of neuron n be denoted A_n^0 , and let the $A_n - A_n^0$ be denoted as δ_n . We define our contributions C_{xy} to satisfy the following properties.

2.1. Summation to δ

For any set of neurons S whose activations are minimally sufficient to compute the activation of y (that is, if we know the activations of S , we can compute the activation of y , and there is no set $S' \subset S$ such that S' is sufficient to compute the activation of y - in layman's terms, S is a full set of non-redundant inputs to y), the following property holds:

$$\sum_{s \in S} C_{sy} = \delta_y \quad (1)$$

That is, the sum over all the contributions of neurons in S to y equals the difference-from-reference of y .

2.2. Linear composition

Let O_x represent the output neurons of x . The following property holds:

$$C_{xy} = \sum_{o \in O_x} \frac{C_{xo}}{\delta_o} C_{oy} \quad (2)$$

In layman's terms, each neuron 'inherits' a contribution through its outputs in proportion to how much that neuron contributes to the difference-from-reference of the output.

2.3. Backpropagation Rules

We show that the contributions as defined above can be computed using the following rules (which can be implemented to run on a GPU). The computation is reminiscent of the chain rule used during gradient backpropagation, as equation 2 makes it possible to start with contribution scores of later layers and use them to find the contribution scores of preceding layers. To avoid issues of numerical stability when δ_n for a particular neuron is small, rather than computing the contribution scores explicitly, we instead compute *multipliers* m_{xy} that, when multiplied with the difference-from-reference, give the contribution:

$$m_{xy} \delta_x = C_{xy} \quad (3)$$

Let t represent the target neuron that we intend to compute contributions to, and let O_x represent the set of outputs of x . We show that:

$$m_{xt} = \sum_{y \in O_x} m_{xy} m_{yt} \quad (4)$$

The equation above follows from the linear composition property and the definition of the multipliers, as proved be-

low:

$$\begin{aligned} C_{xt} &= \sum_{y \in O_x} \frac{C_{xy}}{\delta_y} C_{yt} \\ m_{xt} \delta_x &= \sum_{y \in O_x} \frac{C_{xy}}{\delta_y} (m_{yt} \delta_y) \\ &= \sum_{y \in O_x} C_{xy} m_{yt} \\ m_{xt} &= \sum_{y \in O_x} \frac{C_{xy}}{\delta_x} m_{yt} \\ m_{xt} &= \sum_{y \in O_x} m_{xy} m_{yt} \end{aligned} \quad (5)$$

In the equations below, I_y denotes the set of inputs of y .

2.3.1. AFFINE FUNCTIONS

Let

$$A_y = \left(\sum_{x \in I_y} w_{xy} A_x \right) + b \quad (6)$$

Then

$$m_{xy} = w_{xy} \quad (7)$$

Proof. We show that

$$\delta_y = \sum_{x \in I_y} m_{xy} \delta_x \quad (8)$$

Using the fact that $A_n = A_n^0 + \delta_n$, we have:

$$\begin{aligned} (A_y^0 + \delta_y) &= \left(\sum_{x \in I_y} w_{xy} (A_x^0 + \delta_x) \right) + b \\ &= \left(\sum_{x \in I_y} w_{xy} A_x^0 \right) + b + \sum_{x \in I_y} w_{xy} \delta_x \end{aligned} \quad (9)$$

We also note that the reference activation A_y^0 can be found as follows:

$$A_y^0 = \left(\sum_{x \in I_y} w_{xy} A_x^0 \right) + b \quad (10)$$

Thus, canceling out A_y^0 yields:

$$\begin{aligned} \delta_y &= \sum_{x \in I_y} w_{xy} \delta_x \\ &= \sum_{x \in I_y} m_{xy} \delta_x \end{aligned} \quad (11)$$

2.3.2. MAX POOLING OVER CONVOLUTIONAL FILTERS

We consider the case of maxpooling over a set of neurons in a convolutional layer which all belong to the same convolutional filter:

$$A_y = \max_{x \in I_y} A_x \quad (12)$$

We make the assumption that all instances of a single convolutional filter have the same reference activation, i.e. A_x^0 is the same for all x in I_y (this would not hold if the chosen reference input were not uniform across all positions; the appropriate multiplier in that situation can be computed easily, but is omitted here for brevity). Then we have:

$$m_{xy} = \mathbf{1}\{A_x = A_y\} \quad (13)$$

Where $\mathbf{1}\{\cdot\}$ is the indicator function. If a symbolic computation package is used, then the gradient of y with respect to x can be used in place of $\mathbf{1}\{A_x = A_y\}$.

Proof. Because A_x^0 is assumed to be the same for all x , we have that $A_y^0 = A_x^0$. As $A_y = A_y^0 + \delta_y$, we get:

$$\begin{aligned} \delta_y &= A_y - A_y^0 \\ &= A_y - A_x^0 \\ &= \mathbf{1}\{A_x = A_y\}(A_x - A_x^0) \\ &= \mathbf{1}\{A_x = A_y\}\delta_x \\ &= m_{xy}\delta_x \end{aligned} \quad (14)$$

2.3.3. INDIVIDUAL NONLINEARITY

Here we describe the rule for nonlinear transformations such as the ReLU, sigmoid and tanh which are applied to individual inputs, typically following some affine transformation. We have:

$$A_y = f(A_x) \quad (15)$$

Where f is the nonlinear transformation. When δ_x is large in magnitude, we can compute m_{xy} according to its definition as follows:

$$m_{xy} = \frac{\delta_y}{\delta_x} \quad (16)$$

When δ_x is small, the term on the right approaches $\frac{0}{0}$. However, as $\delta_y \rightarrow 0$ as $\delta_x \rightarrow 0$, we have:

$$\lim_{\delta_x \rightarrow 0} \frac{\delta_y}{\delta_x} = \frac{d\delta_y}{d\delta_x} \quad (17)$$

Thus, when δ_x is close to zero, m_{xy} is simply the derivative of δ_y with respect to δ_x . Because $\delta_n = A_n - A_n^0$ and A_n^0 is a constant, this is the same as $\frac{dA_y}{dA_x}$ when δ_x is near 0.

2.4. A note on final activation layers

Activation functions such as a softmax or a sigmoid have a maximum δ of 1.0. Due to the *summation to δ* property, the

contribution scores for individual features are lower when there are several redundant features present. As an example, consider $A_t = \sigma(A_y)$ (where *sigma* is the sigmoid transformation) and $A_y = A_{x_1} + A_{x_2}$. Let the default activations of the inputs be $A_{x_1}^0 = A_{x_2}^0 = 0$. When $x_1 = 100$ and $x_2 = 0$, we have $C_{x_1 t} = 0.5$. However, when both $x_1 = 100$ and $x_2 = 100$, we have $C_{x_1 t} = C_{x_2 t} = 0.25$. To avoid this attenuation of contribution in the presence of redundant inputs, we can use the contributions to y rather than t ; in both cases, $C_{x_1 y} = 100$.

2.5. A note on Softmax activation

Let $t_1, t_2 \dots t_n$ represent the output of a softmax transformation on the nodes $y_1, y_2 \dots y_n$, such that:

$$A_{t_i} = \frac{e^{A_{y_i}}}{\sum_{i'=1}^n e^{A_{y_{i'}}}} \quad (18)$$

Here, $A_{y_1} \dots A_{y_n}$ are affine functions of their inputs. Let x represent a neuron that is an input to $A_{y_1} \dots A_{y_n}$, and let w_{xy_i} represent the coefficient of A_x in A_{y_i} . Because $A_{y_1} \dots A_{y_n}$ are followed by a softmax transformation, if w_{xy_i} is the same for all y_i (that is, x contributes equally to all y_i), then x effectively has zero contribution to A_{t_i} . This can be observed by substituting $A_{y_i} = w_{xy_i} A_x + r_{y_i}$ in the expression for A_{t_i} and canceling out $e^{w_{xy_i} A_x}$ (here, r_{y_i} is the sum of all the remaining terms in the affine expression for A_{y_i})

$$\begin{aligned} A_{t_i} &= \frac{e^{A_{y_i}}}{\sum_{i'=1}^n e^{A_{y_{i'}}}} \\ &= \frac{e^{w_{xy_i} A_x + r_{y_i}}}{\sum_{i'=1}^n e^{w_{xy_{i'}} A_x + r_{y_{i'}}}} \\ &= \frac{e^{w_{xy_i} A_x + r_{y_i}}}{\sum_{i'=1}^n e^{w_{xy_{i'}} A_x + r_{y_{i'}}}} \\ &= \frac{e^{r_{y_i}}}{\sum_{i'=1}^n e^{r_{y_{i'}}}} \end{aligned} \quad (19)$$

As mentioned in the previous subsection, in order to avoid attenuation of signal for highly confident predictions, we should compute C_{xy_i} rather than C_{xt_i} . One way to ensure that C_{xy_i} is zero if w_{xy_i} is the same for all y_i is to mean-normalized the weights as follows:

$$\bar{w}_{xy_i} = w_{xy_i} - \frac{1}{n} \sum_{i'=1}^n w_{xy_{i'}} \quad (20)$$

This transformation will not affect the output of the softmax, but will ensure that the LIFTPAD scores are zero when a particular node contributes equally to all softmax classes.

2.6. Weight normalization for constrained inputs

Let y be a neuron with some subset of inputs S_y that are constrained such that $\sum_{x \in S_y} A_x = c$ (for example, one-hot encoded input satisfies the constraint $\sum_{x \in S_y} A_x = 1$, and a convolutional neuron operating on one-hot encoded rows has one constraint per column that it sees). Let the weights from x to y be denoted w_{xy} and let b_y be the bias of y . It is advisable to use normalized weights $\bar{w}_{xy} = w_{xy} - \mu$ and bias $\bar{b}_y = b_y + c\mu$, where μ is the mean over all w_{xy} . We note that this maintains the output of the neural net because, for any constant μ :

$$\begin{aligned} A_y &= \left(\sum A_x (\bar{w}_{xy} - \mu) \right) + (b_y + c\mu) \\ &= \left(\sum A_x w_{xy} \right) - \left(\sum A_x \mu \right) + (b_y + c\mu) \\ &= \left(\sum A_x w_{xy} \right) - c\mu + (b_y + c\mu) \\ &= \left(\sum A_x w_{xy} \right) + b_y \end{aligned} \quad (21)$$

The normalization is desirable because, for affine functions, the multipliers m_{xy} are equal to the weights w_{xy} and are thus sensitive to μ . To take the example of a convolutional neuron operating on one-hot encoded rows: by mean-normalizing w_{xy} for each column in the filter, one can ensure that the contributions C_{xy} from some columns are not systematically overestimated or underestimated relative to the contributions from other columns.

3. Results

3.1. Tiny Imagenet

A model with the VGG16 [cite] architecture was trained on a scaled-down version of the Imagenet dataset, dubbed ‘Tiny Imagenet’. The images were 64×64 in dimension and belonged to one of 200 output classes. Pictured are importance scores computed using three different approaches.

3.2. Genomics

We apply LIFTPAD to models that classify genomic sequences. The positive class requires that two particular DNA patterns - ‘GATA’ and ‘CAGATG’ - appear in the sequence together. They may appear at any spacing. The negative class has just one of the two patterns appearing either once or twice, also at any spacing. We simulate 20,000 sequences of length 200 for the positive set and 40,000 sequences for the negative set split equally between the two patterns. Outside the core patterns we randomly sample the remaining sequence with equal frequency for the four bases A, C, G and T. We trained a model on these sequences with 20 convolutional filters of length 15 with a PReLU and a max pool layer of stride 50 followed by two fully connected layers of size 200. This model performs well with auROC

of 0.907. The misclassified examples primarily occur when one of the patterns erroneously arises in the randomly sampled background.

We then run LIFTPAD to assign an importance score to each base in the correctly predicted sequences. We also use a gradient-based method for comparison. In Figure 1 we illustrate that LIFTPAD clearly discovers both patterns and assigns them large importance scores as compared to the surrounding sequence. The gradient-based method discovers only one of the two important patterns despite a correct prediction.

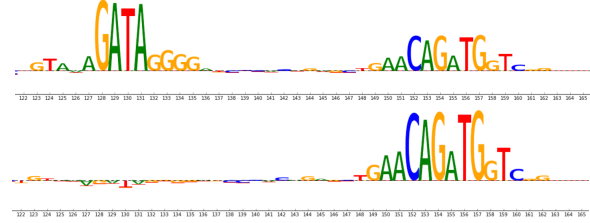


Figure 1. LIFTPAD scores (top) and gradient-based scores (bottom) are plotted for each position in the DNA sequence and colored by the DNA base.

3.3. Extension to Recurrent Neural Networks

LIFTPAD is especially appropriate for computing importance scores in Recurrent architectures that involve gated memory units because, unlike with gradients, LIFTPAD contribution scores do not vanish at the extremes of a sigmoid or tanh function. To find importance scores, the computation would have to be done on the unrolled RNN graph.

3.4. Author contributions

AS and PG conceived of LIFTPAD. AS implemented LIFTPAD in software. PG led the application to genomics. AYS led the application to Tiny Imagenet. AK provided guidance and feedback. AS, PG, AYS and AK prepared the manuscript.

4. References