

---

# LEARNING HOW TO EXPLAIN NEURAL NETWORKS: PATTERNNET AND PATTERNATTRIBUTION

**Pieter-Jan Kindermans\***

Google Brain

pikinder@google.com

**Kristof T. Schütt & Maximilian Alber**

TU Berlin

{maximilian.alber,kristof.schuett}@tu-berlin.de

**Klaus-Robert Müller<sup>†</sup>**

TU Berlin

klaus-robert.mueller@tu-berlin.de

**Dumitru Erhan & Been Kim**

Google Brain

{dumitru,beenkim}@google.com

**Sven Dähne<sup>‡</sup>**

TU Berlin

sven.daehne@tu-berlin.de

## ABSTRACT

DeConvNet, Guided BackProp, LRP, were invented to better understand deep neural networks. We show that these methods do not produce the theoretically correct explanation for a linear model. Yet they are used on multi-layer networks with millions of parameters. This is a cause for concern since linear models are simple neural networks. We argue that explanation methods for neural nets should work reliably in the limit of simplicity, the linear models. Based on our analysis of linear models we propose a generalization that yields two explanation techniques (PatternNet and PatternAttribution) that are theoretically sound for linear models and produce improved explanations for deep networks.

## 1 INTRODUCTION

Deep learning made a huge impact on a wide variety of applications (LeCun et al., 2015; Schmidhuber, 2015; Krizhevsky et al., 2012; Mnih et al., 2015; Silver et al., 2016; Sutskever et al., 2014) and recent neural network classifiers have become extremely good at detecting relevant *signals* (say, the presence of a cat) contained in input data points such as images by filtering out all other, non-relevant and *distracting* components also present in the data. This separation of signal and distractors is achieved by passing the input through many layers with millions of parameters and nonlinear activation functions in between until finally at the output layer, these models yield a highly condensed version of the signal, e.g. a single number indicating the probability of a cat being in the image.

While deep neural networks learn efficient and powerful representations, they are often considered a ‘black-box’. In order to better understand classifier decisions and to gain insight into how these models operate, a variety techniques have been proposed (Simonyan et al., 2014; Yosinski et al., 2015; Nguyen et al., 2016; Baehrens et al., 2010; Bach et al., 2015; Montavon et al., 2017; Zeiler & Fergus, 2014; Springenberg et al., 2015; Zintgraf et al., 2017; Sundararajan et al., 2017; Smilkov et al., 2017).

The aforementioned methods for explaining classifier decisions operate under the assumption that it is possible to propagate the condensed output signal back through the classifier to arrive at something that shows how the relevant signal was encoded in the input and thereby explains the classifier decision. Simply put, if the classifier detected a cat, the visualization should point to the cat-relevant aspects of the input image from the perspective of the network. Techniques that are based on this

---

\*Part of this work was done at TU Berlin, part of the work was part of the Google Brain Residency program.

<sup>†</sup>KRM is also with Korea University and Max Planck Institute for Informatics, Saarbrücken, Germany

<sup>‡</sup>Sven Dähne is now at Amazon

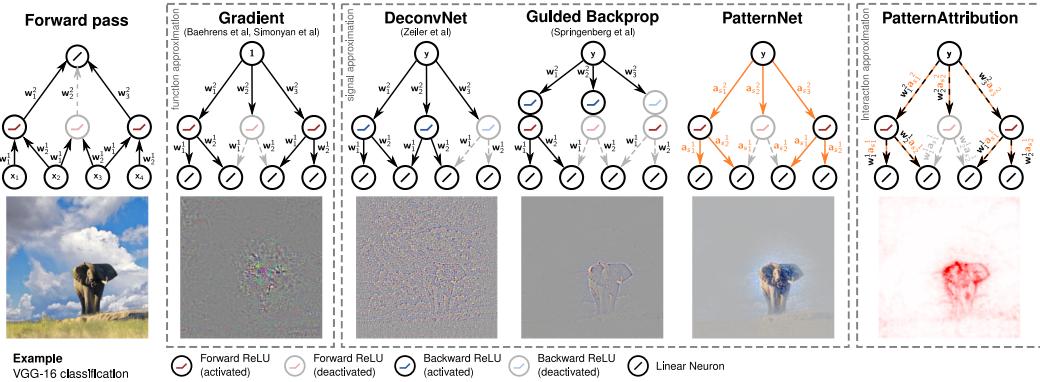


Figure 1: Illustration of explanation approaches. Function and signal approximators visualize the explanation using the original color channels. The attribution is visualized as a heat map of pixel-wise contributions to the output

principle include saliency maps from network gradients (Baehrens et al., 2010; Simonyan et al., 2014), DeConvNet (Zeiler & Fergus, 2014, DCN), Guided BackProp (Springenberg et al., 2015, GBP), Layer-wise Relevance Propagation (Bach et al., 2015, LRP) and the Deep Taylor Decomposition (Montavon et al., 2017, DTD), Integrated Gradients (Sundararajan et al., 2017) and SmoothGrad (Smilkov et al., 2017).

The merit of explanation methods is often proven by applying them to state-of-the-art deep learning models in the context of high dimensional real world data, such as ImageNet. Here we begin with a different approach. We first take a step back and analyze explanation methods in the context of the simplest neural network setting: a purely linear model and data stemming from a linear generative model. We chose this simplified setup because it allows us to (i) fully control how signal and distractor components are encoded in the input data and (ii) analytically track how the resulting explanation relates to the known signal component. This analysis allows us to highlight shortcomings of current explanation approaches that carry over to non-linear models as well.

On the basis of our findings, we then propose PatternNet and PatternAttribution, which alleviate these flaws. Finally we apply our methods to practically relevant networks and datasets, and show that our approach produces qualitatively improved signal visualizations and attributions (see Fig. 2 and Fig. 4b). In addition to the qualitative evaluation, we also experimentally verify whether our proposed theoretical model holds up empirically (see Fig. 3).

The remainder of the paper is structured as follows: visualization in linear models is analyzed in section 2. Section 3 relates this analysis to existing approaches for neural network visualization. Section 4 introduces PatternNet and PatternAttribution, which then is evaluated in section 5 before concluding.

**Notation and scope** Scalars are lowercase letters ( $i$ ), column vectors are bold ( $\mathbf{u}$ ), element-wise multiplication is ( $\odot$ ). The covariance between  $\mathbf{u}$  and  $\mathbf{v}$  is  $\text{cov}[\mathbf{u}, \mathbf{v}]$ , the covariance of  $\mathbf{u}$  and  $i$  is  $\text{cov}[\mathbf{u}, i]$ . The variance of a scalar random variable  $i$  is  $\sigma_i^2$ . Estimates of random variables will have a hat ( $\hat{u}$ ). We analyze neural networks excluding the final soft-max output layer. To allow for analytical treatment, we only consider networks with linear neurons optionally followed by a rectified linear unit (ReLU), max-pooling or soft-max. We analyze linear neurons and nonlinearities independently such that every neuron has its own weight vector. These restrictions are similar to those in the saliency map (Simonyan et al., 2014), DCN (Zeiler & Fergus, 2014), GBP (Springenberg et al., 2015), LRP (Bach et al., 2015) and DTD (Montavon et al., 2017). Without loss of generality, biases are considered constant neurons to enhance clarity.

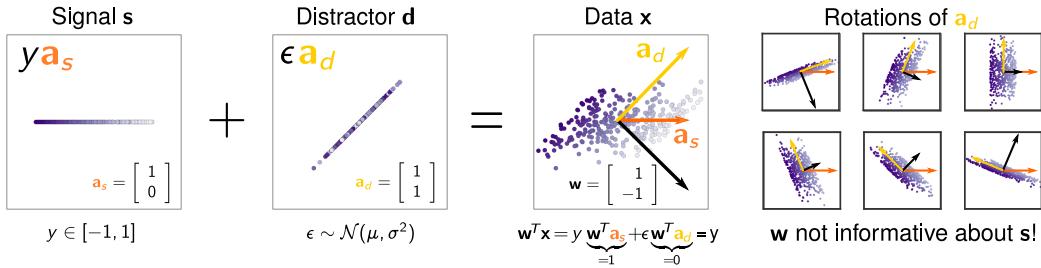


Figure 2: For linear models, i.e., a simple neural network, the weight vector does **not** explain the signal it detects Haufe et al. (2014). The data  $\mathbf{x} = y\mathbf{a}_s + \epsilon\mathbf{a}_d$  is color-coded w.r.t. the output  $y = \mathbf{w}^T \mathbf{x}$ . Only the signal  $s = y\mathbf{a}_s$  contributes to  $y$ . The weight vector  $\mathbf{w}$  does not agree with the signal direction, since its primary objective is canceling the distractor. Therefore, rotations of the basis vector  $\mathbf{a}_d$  of the distractor with constant signal  $s$  lead to rotations of the weight vector (**right**).  $\mathbf{w}$  not informative about  $s$ !

## 2 UNDERSTANDING LINEAR MODELS

Before moving to deep networks, we analyze the behavior of a linear model (see Fig. 2). Consider the following toy example where we generate data  $\mathbf{x}$  as:

$$\begin{aligned} \mathbf{x} &= \mathbf{s} + \mathbf{d} & \mathbf{s} &= \mathbf{a}_s y, & \text{with } \mathbf{a}_s &= (1, 0)^T, & y &\in [-1, 1] \\ \mathbf{d} &= \mathbf{a}_d \epsilon, & \text{with } \mathbf{a}_d &= (1, 1)^T, & \epsilon &\sim \mathcal{N}(\mu, \sigma^2). \end{aligned}$$

We train a linear regression model to extract  $y$  from  $\mathbf{x}$ . By construction,  $\mathbf{s}$  is the *signal* in our data, i.e., the part of  $\mathbf{x}$  containing information about  $y$ . Using the terminology of Haufe et al. (2014) the *distractor*  $\mathbf{d}$  obfuscates the signal making the detection task more difficult. To optimally extract  $y$ , our model has to be able to filter out the distractor  $\mathbf{d}$ . This is why the weight vector is also called the *filter*. In the example,  $\mathbf{w} = [1, -1]^T$  fulfills this task.

From this simple example, we can make several observations: The optimal weight vector  $\mathbf{w}$  does *not* align, in general, with the signal direction  $\mathbf{a}_s$ , but tries to filter the contribution of the distractor (see Fig. 2). This optimally solved when the weight vector is orthogonal to the distractor  $\mathbf{w}^T \mathbf{d} = 0$ . Therefore, when the direction of the distractor  $\mathbf{a}_d$  changes,  $\mathbf{w}$  must follow, as illustrated on the right hand side of the figure. On the other hand, a change in signal direction  $\mathbf{a}_s$  can be compensated for by a change in sign and magnitude of  $\mathbf{w}$  such that  $\mathbf{w}^T \mathbf{a}_s = 1$ , but the direction stays constant.

This implies that in the situation where we have a signal and distractor(s), the direction of the weight vector in a linear model is largely determined by the distractor. This is essential to understand how linear models operate. This also indicates that given only the weight vector, we cannot know what part of the input produces the output  $y$ . This is the direction  $\mathbf{a}_s$  and must be learned from data. Please note that the linear problem above is convex, therefore a weight vector obtained by optimizing the model would converge to the analytical solution defined above.

Now assume that we have no distractor but instead we have additive isotropic Gaussian noise. It is easy to verify that the mean of the noise is of little importance since it can be compensated for with a bias change. Therefore, we only have to consider case where the noise is zero mean. Because isotropic Gaussian noise does not contain any correlations or structure, the only way to cancel it out is by averaging over different measurements. It is *not* possible to cancel it out effectively by using a well-chosen weight vector. However, it is well known that adding Gaussian noise shrinks the weight vector and corresponds to L2 regularization. In the absence of a structured distractor, the smallest weight vector  $\mathbf{w}$  such that  $\mathbf{w}^T \mathbf{a}_s = 1$  is the one in the direction of the signal. Therefore in practice both these effects influence the actual weight vector.

Considering the reasoning above, we have to wonder under which conditions we are working in a deep neural network. Especially since DeConvNet and Guided BackProp produce crisp visualizations using (modified) gradients. For this reason we will perform the following quantitative and qualitative experiments which indicate that our theory also holds for a deep network:

- 
- In Fig. 3 we have evaluated how well the weight vector or a learned direction captures the information content in the input of every single neuron in VGG16. This experiment empirically shows that a learned direction captures more information than the direction defined by the weight vector. This indicates that we are working (largely) in the distractor-regime.
  - This experiment is confirmed by an image degradation experiment in Fig. 4a.
  - It is also corroborated by the qualitative inspection of the visualizations in Fig. 1, Fig. 4b and Fig. 5.

Finally, there is also an intuitive argument. Neural networks are considered layer-wise feature extractors that add more invariances as we move through the layers. Since cancelling out a distractor is adding an invariance, the proposed theory fits this interpretation well.

Before moving on to the discussion of interpretability methods, we would like to remind the reader of the terminology that is used throughout this manuscript: The *filter  $w$*  tells us how to extract the output  $y$  optimally from data  $x$ . The *pattern  $a_s$*  is the direction in the data along which the desired output  $y$  varies. Both constitute the *signal  $s = a_s y$* , i.e., the contributing part of  $x$ . The *distractor  $d$*  is the component of the data that does not contain information about the desired output.

### 3 OVERVIEW OF EXPLANATION APPROACHES AND THEIR BEHAVIOR

In this section, we take a look at a subset of explanation methods for individual classifier decisions and discuss how they are connected to our analysis of linear models in the previous section. Fig. 1 gives an overview of the different types of explanation methods which can be divided into function, signal and attribution visualizations. These three groups all present different information about the network and complement each other.

**Functions – gradients, saliency map** Explaining the function in input space corresponds to describing the operations the model uses to extract  $y$  from  $x$ . Since deep neural networks are highly nonlinear, this can only be approximated. The saliency map estimates how moving along a particular direction in input space influences  $y$  (i.e., sensitivity analysis) where the direction is given by the model gradient (Baehrens et al., 2010; Simonyan et al., 2014). In case of a linear model  $y = w^T x$ , the saliency map reduces to analyzing the weights  $\partial y / \partial x = w$ . Since it is mostly determined by the distractor, as demonstrated above, it is not representing the signal. It tells us how to extract the signal, not what the signal is in a deep neural network.

**Signal – DeConvNet, Guided BackProp, PatternNet** The signal  $s$  detected by the neural network is the component of the data that caused the networks activations. Zeiler & Fergus (2014) formulated the goal of these methods as "[...] to map these activities back to the input pixel space, showing what input pattern originally caused a given activation in the feature maps".

In a linear model, the signal corresponds to  $s = a_s y$ . The pattern  $a_s$  contains the signal direction, i.e., it tells us where a change of the output variable is expected to be measurable in the input (Haufe et al., 2014). Attempts to visualize the signal for deep neural networks were made using DeConvNet (Zeiler & Fergus, 2014) and Guided BackProp (Springenberg et al., 2015). These use the same algorithm as the saliency map, but treat the rectifiers differently (see Fig. 1): DeConvNet leaves out the rectifiers from the forward pass, but adds additional ReLUs after each deconvolution, while Guided BackProp uses the ReLUs from the forward pass as well as additional ones. The back-projections for the linear components of the network correspond to a superposition of what are assumed to be the signal directions of each neuron. For this reason, these projections must be seen as an approximation of the features that activated the higher layer neuron. It is not a reconstruction in input space (Zeiler & Fergus, 2014).

For the simplest of neural networks – the linear model – these visualizations reduce to the gradient<sup>1</sup>. They show the filter  $w$  and *neither* the pattern  $a_s$ , *nor* the signal  $s$ . Hence, DeConvNet and Guided BackProp do not guarantee to produce the detected signal for a linear model, which is proven by our

---

<sup>1</sup>In tensorflow terminoloy: linear model on MNIST can be seen as a convolutional neural network with VALID padding and a 28 by 28 filter size.

---

toy example in Fig. 2. Since they do produce compelling visualizations, we will later investigate whether the direction of the filter  $\mathbf{w}$  coincides with the direction of the signal  $\mathbf{s}$ . We will show that this is *not* the case and propose a new approach, PatternNet (see Fig. 1), to estimate the correct direction that improves upon the DeConvNet and Guided BackProp visualizations.

**Attribution – LRP, Deep Taylor Decomposition, PatternAttribution** Finally, we can look at how much the signal dimensions contribute to the output through the layers. This will be referred to as the *attribution*. For a linear model, the optimal attribution would be obtained by element-wise multiplying the signal with the weight vector:  $\mathbf{r}^{input} = \mathbf{w} \odot \mathbf{ay}$ , with  $\odot$  the element-wise multiplication. Bach et al. (2015) introduced *layer-wise relevance propagation* (LRP) as a decomposition of pixel-wise contributions (called *relevances*). Montavon et al. (2017) extended this idea and proposed the deep Taylor decomposition (DTD). The key idea of DTD is to decompose the activation of a neuron in terms of contributions from its inputs. This is achieved using a first-order Taylor expansion around a root point  $\mathbf{x}_0$  with  $\mathbf{w}^T \mathbf{x}_0 = 0$ . The relevance of the selected output neuron  $i$  is initialized with its output from the forward pass. The relevance from neuron  $i$  in layer  $l$  is re-distributed towards its input as:

$$r_i^{output} = y, \quad r_{j \neq i}^{output} = 0, \quad r_i^{l-1} = \frac{\mathbf{w} \odot (\mathbf{x} - \mathbf{x}_0)}{\mathbf{w}^T \mathbf{x}} r_i^l.$$

Here we can safely assume that  $\mathbf{w}^T \mathbf{x} > 0$  because a non-active ReLU unit from the forward pass stops the re-distribution in the backward pass. This is identical to how a ReLU stops the propagation of the gradient. The difficulty in the application of the deep Taylor decomposition is the choice of the root point  $\mathbf{x}_0$ , for which many options are available. It is important to recognize at this point that selecting a root point for the DTD corresponds to estimating the distractor  $\mathbf{x}_0 = \mathbf{d}$  and, by that, the signal  $\hat{\mathbf{s}} = \mathbf{x} - \mathbf{x}_0$ . PatternAttribution is a DTD extension that learns from data how to set the root point.

Summarizing, the **function** extracts the **signal** from the data by removing the distractor. The **attribute** of output values to input dimensions shows how much an individual component of the signal contributes to the output, which is what LRP calls *relevance*.

## 4 LEARNING TO ESTIMATE THE SIGNAL

Visualizing the function has proven to be straightforward (Baehrens et al., 2010; Simonyan et al., 2014). In contrast, visualizing the signal (Haufe et al., 2014; Zeiler & Fergus, 2014; Springenberg et al., 2015) and the attribution (Bach et al., 2015; Montavon et al., 2017; Sundararajan et al., 2017) is more difficult. It requires a good estimate of what is the signal and what is the distractor. In the following section we first propose a quality measure for neuron-wise signal estimators. This allows us to evaluate existing approaches and, finally, derive signal estimators that optimize this criterion. These estimators will then be used to explain the signal (PatternNet) and the attribution (PatternAttribution). All mentioned techniques as well as our proposed signal estimators treat neurons independently, i.e., the full explanation will be a superposition of neuron-wise explanations.

### 4.1 QUALITY CRITERION FOR SIGNAL ESTIMATORS

Recall that the input data  $\mathbf{x}$  comprises both signal and distractor:  $\mathbf{x} = \mathbf{s} + \mathbf{d}$ , and that the signal contributes to the output but the distractor does not. Assuming the filter  $\mathbf{w}$  has been trained sufficiently well to extract  $y$ , we have

$$\mathbf{w}^T \mathbf{x} = y, \quad \mathbf{w}^T \mathbf{s} = y, \quad \mathbf{w}^T \mathbf{d} = 0.$$

Note that estimating the signal based on these conditions alone is an ill-posed problem. We could limit ourselves to linear estimators of the form  $\hat{\mathbf{s}} = \mathbf{u}(\mathbf{w}^T \mathbf{u})^{-1} y$ , with  $\mathbf{u}$  a random vector such that  $\mathbf{w}^T \mathbf{u} \neq 0$ . For such an estimator, the signal estimate  $\hat{\mathbf{s}} = \mathbf{u}(\mathbf{w}^T \mathbf{u})^{-1} y$  satisfies  $\mathbf{w}^T \hat{\mathbf{s}} = y$ . This implies the existence of an infinite number of possible rules for the DTD as well as infinitely many back-projections for the DeConvNet family.

To alleviate this issue, we introduce the following quality measure  $\rho$  for a signal estimator  $S(\mathbf{x}) = \hat{\mathbf{s}}$  that will be written with explicit variances and covariances using the shorthands  $\hat{\mathbf{d}} = \mathbf{x} - S(\mathbf{x})$  and

---


$$y = \mathbf{w}^T \mathbf{x}$$

$$\rho(S) = 1 - \max_{\mathbf{v}} \text{corr}(\mathbf{w}^T \mathbf{x}, \mathbf{v}^T (\mathbf{x} - S(\mathbf{x}))) = 1 - \max_{\mathbf{v}} \frac{\mathbf{v}^T \text{cov}[\hat{\mathbf{d}}, y]}{\sqrt{\sigma_{\mathbf{v}^T \hat{\mathbf{d}}}^2 \sigma_y^2}}. \quad (1)$$

This criterion introduces an additional constraint by measuring how much information about  $y$  can be reconstructed from the residuals  $\mathbf{x} - \hat{\mathbf{s}}$  using a linear projection. The best signal estimators remove most of the information in the residuals and thus yield large  $\rho(S)$ . Since the correlation is invariant to scaling, we constrain  $\mathbf{v}^T \hat{\mathbf{d}}$  to have variance  $\sigma_{\mathbf{v}^T \hat{\mathbf{d}}}^2 = \sigma_y^2$ . Finding the optimal  $\mathbf{v}$  for a fixed  $S(\mathbf{x})$  amounts to a least-squares regression from  $\hat{\mathbf{d}}$  to  $y$ . This enables us to assess the quality of signal estimators efficiently.

## 4.2 EXISTING SIGNAL ESTIMATORS

Let us now discuss two signal estimators that have been used in previous approaches.

**$S_x$  – the identity estimator** The naive approach to signal estimation is to assume the entire data is signal and there are no distractors:

$$S_x(\mathbf{x}) = \mathbf{x}.$$

With this being plugged into the deep Taylor framework, we obtain the  $z$ -rule (Montavon et al., 2017) which is equivalent to LRP (Bach et al., 2015). For a linear model, this corresponds to  $\mathbf{r} = \mathbf{w} \odot \mathbf{x}$  as the attribution. It can be shown that for ReLU and max-pooling networks, the  $z$ -rule reduces to the element-wise multiplication of the input and the saliency map (Shrikumar et al., 2016; Kindermans et al., 2016). This means that for a whole network, the assumed signal is simply the original input image. It also implies that, if there are distractors present in the data, they are included in the attribution:

$$\mathbf{r} = \mathbf{w} \odot \mathbf{x} = \mathbf{w} \odot \mathbf{s} + \mathbf{w} \odot \mathbf{d}.$$

When moving through the layers by applying the filters  $\mathbf{w}$  during the forward pass, the contributions from the distractor  $\mathbf{d}$  are cancelled out. However, they cannot be cancelled in the backward pass by the element-wise multiplication. The distractor contributions  $\mathbf{w} \odot \mathbf{d}$  that are included in the LRP explanation cause the noisy nature of the visualizations based on the  $z$ -rule.

**$S_w$  – the filter based estimator** The implicit assumption made by DeConvNet and Guided Back-Prop is that the detected signal varies in the direction of the weight vector  $\mathbf{w}$ . This weight vector has to be normalized in order to be a valid signal estimator. In the deep Taylor decomposition framework this corresponds to the  $\mathbf{w}^2$ -rule and results in the following signal estimator:

$$S_w(\mathbf{x}) = \frac{\mathbf{w}}{\mathbf{w}^T \mathbf{w}} \mathbf{w}^T \mathbf{x}.$$

For a linear model, this produces an attribution of the form  $\frac{\mathbf{w} \odot \mathbf{w}}{\mathbf{w}^T \mathbf{w}} y$ . This estimator does not reconstruct the proper signal in the toy example of section 2. Empirically it is also sub-optimal in our experiment in Fig. 3.

## 4.3 PATTERNNET AND PATTERNATTRIBUTION

We suggest to learn the signal estimator  $S$  from data by optimizing the previously established criterion. A signal estimator  $S$  is optimal with respect to Eq. (1) if the correlation is zero for all possible  $\mathbf{v}$ :  $\forall \mathbf{v}, \text{cov}[y, \hat{\mathbf{d}}] \mathbf{v} = \mathbf{0}$ . This is the case when there is no covariance between  $y$  and  $\hat{\mathbf{d}}$ . Because of linearity of the covariance and since  $\hat{\mathbf{d}} = \mathbf{x} - S(\mathbf{x})$  the above condition leads to

$$\text{cov}[y, \hat{\mathbf{d}}] = \mathbf{0} \Rightarrow \text{cov}[\mathbf{x}, y] = \text{cov}[S(\mathbf{x}), y]. \quad (2)$$

It is important to recognize that the covariance is a summarizing statistic and consequently the problem can still be solved in multiple ways. We will present two possible solutions to this problem. Note that when optimizing the estimator, the contribution from the bias neuron will be considered 0 since it does not covary with the output  $y$ .

---

**$S_a$  – The linear estimator** A linear neuron can only extract linear signals  $s$  from its input  $x$ . Therefore, we could assume a linear dependency between  $s$  and  $y$ , yielding a signal estimator:

$$S_a(x) = \mathbf{a} \mathbf{w}^T x. \quad (3)$$

Plugging this into Eq. (2) and optimising for  $\mathbf{a}$  yields

$$\text{cov}[x, y] = \text{cov}[\mathbf{a} \mathbf{w}^T x, y] = \mathbf{a} \text{cov}[y, y] \Rightarrow \mathbf{a} = \frac{\text{cov}[x, y]}{\sigma_y^2}.$$

Note that this solution is equivalent to the approach commonly used in neuro-imaging (Haufe et al., 2014) despite different derivation. With this approach we can recover the signal of our toy example in section 2. It is equivalent to the filter-based approach only if the distractors are orthogonal to the signal. We found that the linear estimator works well for the convolutional layers. However, when using this signal estimator with ReLUs in the dense layers, there is still a considerable correlation left in the distractor component (see Fig. 3).

**$S_{a+-}$  – The two-component estimator** To move beyond the linear signal estimator, it is crucial to understand how the rectifier influences the training. Since the gate of the ReLU closes for negative activations, the weights only need to filter the distractor component of neurons with  $y > 0$ . Since this allows the neural network to apply filters locally, we cannot assume a global distractor component. We rather need to distinguish between the positive and negative regime:

$$x = \begin{cases} s_+ + d_+ & \text{if } y > 0 \\ s_- + d_- & \text{otherwise} \end{cases}$$

Even though signal and distractor of the negative regime are canceled by the following ReLU, we still need to make this distinction in order to approximate the signal. Otherwise, information about whether a neuron fired would be retained in the distractor. Thus, we propose the two-component signal estimator:

$$S_{a+-}(x) = \begin{cases} \mathbf{a}_+ \mathbf{w}^T x, & \text{if } \mathbf{w}^T x > 0 \\ \mathbf{a}_- \mathbf{w}^T x, & \text{otherwise} \end{cases} \quad (4)$$

Next, we derive expressions for the patterns  $\mathbf{a}_+$  and  $\mathbf{a}_-$ . We denote expectations over  $x$  within the positive and negative regime with  $\mathbb{E}_+[x]$  and  $\mathbb{E}_-[x]$ , respectively. Let  $\pi_+$  be the expected ratio of inputs  $x$  with  $\mathbf{w}^T x > 0$ . The covariance of data/signal and output become:

$$\text{cov}[x, y] = \pi_+(\mathbb{E}_+[xy] - \mathbb{E}_+[x]\mathbb{E}[y]) + (1 - \pi_+)(\mathbb{E}_-[xy] - \mathbb{E}_-[x]\mathbb{E}[y]) \quad (5)$$

$$\text{cov}[s, y] = \pi_+(\mathbb{E}_+[sy] - \mathbb{E}_+[s]\mathbb{E}[y]) + (1 - \pi_+)(\mathbb{E}_-[sy] - \mathbb{E}_-[s]\mathbb{E}[y]) \quad (6)$$

Assuming both covariances are equal, we can treat the positive and negative regime separately using Eq. (2) to optimize the signal estimator:

$$\mathbb{E}_+[xy] - \mathbb{E}_+[x]\mathbb{E}[y] = \mathbb{E}_+[sy] - \mathbb{E}_+[s]\mathbb{E}[y]$$

Plugging in Eq. (4) and solving for  $\mathbf{a}_+$  yields the required parameter ( $\mathbf{a}_-$  analogous).

$$\mathbf{a}_+ = \frac{\mathbb{E}_+[xy] - \mathbb{E}_+[x]\mathbb{E}[y]}{\mathbf{w}^T \mathbb{E}_+[xy] - \mathbf{w}^T \mathbb{E}_+[x]\mathbb{E}[y]} \quad (7)$$

The solution for  $S_{a+-}$  reduces to the linear estimator when the relation between input and output is linear. Therefore, it solves our introductory linear example correctly.

**PatternNet and PatternAttribution** Based on the presented analysis, we propose PatternNet and PatternAttribution as illustrated in Fig. 1. *PatternNet* yields a layer-wise back-projection of the estimated signal to input space. The signal estimator is approximated as a superposition of neuron-wise, nonlinear signal estimators  $S_{a+-}$  in each layer. It is equal to the computation of the gradient where during the backward pass the weights of the network are replaced by the informative directions. In Fig. 1, a visual improvement over DeConvNet and Guided Backprop is apparent.

*PatternAttribution* exposes the attribution  $\mathbf{w} \odot \mathbf{a}_+$  and improves upon the layer-wise relevance propagation (LRP) framework (Bach et al., 2015). It can be seen as a root point estimator for the Deep-Taylor Decomposition (DTD). Here, the explanation consists of neuron-wise contributions of the estimated *signal* to the classification score. By ignoring the distractor, PatternAttribution can reduce the noise and produces much clearer heat maps. By working out the back-projection steps in the Deep-Taylor Decomposition with the proposed root point selection method, it becomes obvious that PatternAttribution is also analogous to the backpropagation operation. In this case, the weights are replaced during the backward pass by  $\mathbf{w} \odot \mathbf{a}_+$ .

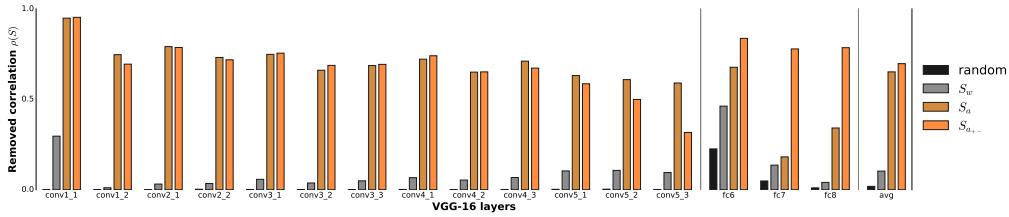
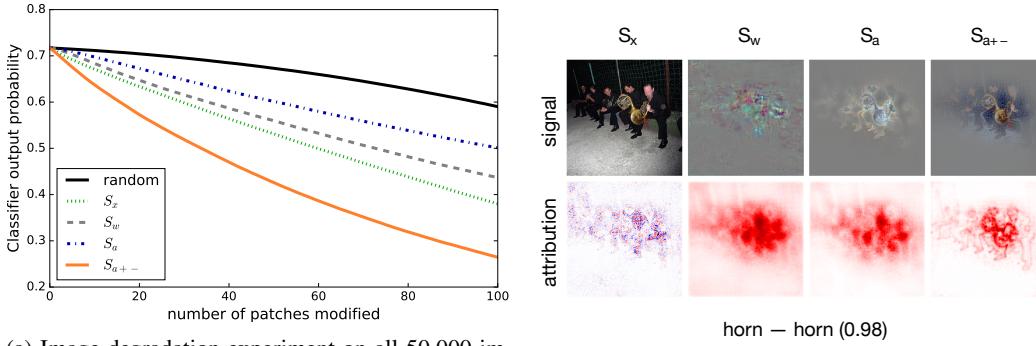


Figure 3: Evaluating  $\rho(S)$  for VGG-16 on ImageNet. Higher values are better. The gradient ( $S_w$ ), linear estimator ( $S_a$ ) and nonlinear estimator ( $S_{a+-}$ ) are compared. An estimator using random directions is the baseline. The network has 5 blocks with 2/3 convolutional layers and 1 max-pooling layer each, followed by 3 dense layers.



(a) Image degradation experiment on all 50.000 images in the ImageNet validation set. The effect on the classifier output is measured. A steeper decrease is better.

(b) Top: signal. Bottom: attribution. For the trivial estimator  $S_x$  the original input is the signal. This is not informative w.r.t. how the network operates.

## 5 EXPERIMENTS AND DISCUSSION

To evaluate the quality of the explanations, we focus on the task of image classification. Nevertheless, our method is not restricted to networks operating on image inputs. We used Theano (Bergstra et al., 2010) and Lasagne (Dieleman et al., 2015) for our implementation. We restrict the analysis to the well-known ImageNet dataset (Russakovsky et al., 2015) using the pre-trained VGG-16 model (Simonyan & Zisserman, 2015). Images were rescaled and cropped to 224x224 pixels. The signal estimators are trained on the first half of the training dataset.

The vector  $v$ , used to measure the quality of the signal estimator  $\rho(x)$  in Eq. (1), is optimized on the second half of the training dataset. This enables us to test the signal estimators for generalization. All the results presented here were obtained using the official validation set of 50000 samples. The validation set was not used for training the signal estimators, nor for training the vector  $v$  to measure the quality. Consequently our results are obtained on previously unseen data.

The linear and the two component signal estimators are obtained by solving their respective closed form solutions (Eq. (4) and Eq. (7)). With a highly parallelized implementation using 4 GPUs this could be done in 3-4 hours. This can be considered reasonable given that several days are required to train the actual network. The quality of a signal estimator is assessed with Eq. (1). Solving it with the closed form solution is computationally prohibitive since it must be repeated for every single weight vector in the network. Therefore we optimize the equivalent least-squares problem using stochastic mini-batch gradient descent with ADAM Kingma & Ba (2015) until convergence. This was implemented on a NVIDIA Tesla K40 and took about 24 hours per optimized signal estimator.

After learning to explain, individual explanations are computationally cheap since they can be implemented as a back-propagation pass with a modified weight vector. As a result, our method produces explanations at least as fast as the work by Dabkowski & Gal (2017) on real time saliency. However, our method has the advantage that it is not only applicable to image models but is a generalization of the theory commonly used in neuroimaging Haufe et al. (2014).

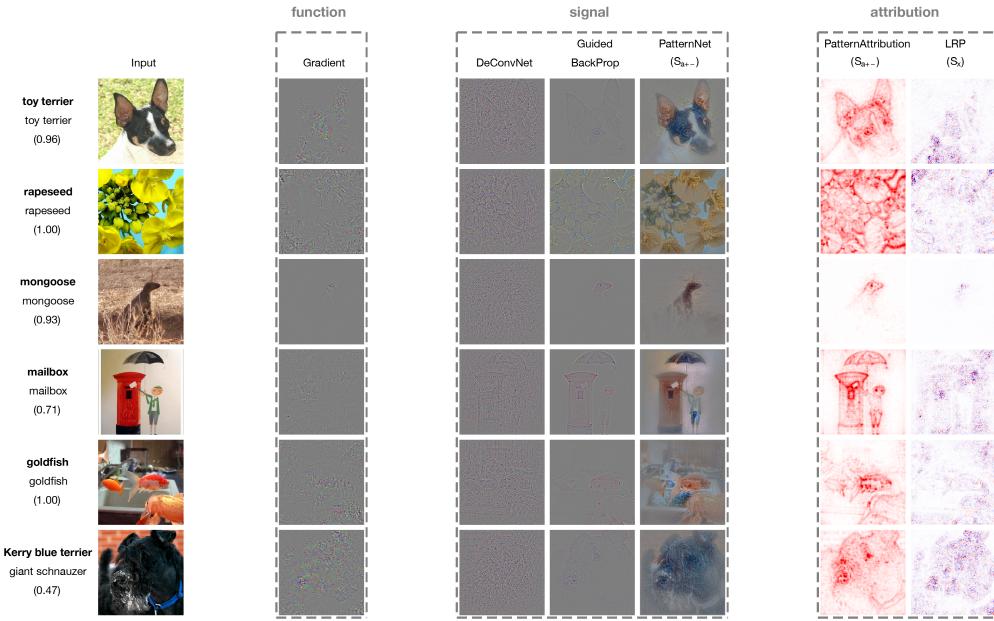


Figure 5: Visualization of random images from ImageNet (validation set). In the leftmost shows column the ground truth, the predicted label and the classifier’s confidence. Methods should only be compared within their group. PatternNet, Guided Backprop, DeConvNet and the Gradient (saliency map) are back-projections to input space with the original color channels. They are normalized using  $x_{norm} = \frac{x}{2 \max|x|} + \frac{1}{2}$  to maximize contrast. LRP and PatternAttribution are heat maps showing pixel-wise contributions. Best viewed in electronic format (zoomed in). The supplementary contains more samples.

**Measuring the quality of signal estimators** In Fig. 3 we present the results from the correlation measure  $\rho(x)$ , where higher values are better. We use random directions as baseline signal estimators. Clearly, this approach removes almost no correlation. The filter-based estimator  $S_w$  succeeds in removing some of the information in the first layer. This indicates that the filters are similar to the patterns in this layer. However, the gradient removes much less information in the higher layers. Overall, it does not perform much better than the random estimator. *This implies that the weights do not correspond to the detected stimulus in a neural network.* Hence the implicit assumptions about the signal made by DeConvNet and Guided BackProp is not valid. The optimized estimators remove much more of the correlations across the board. For convolutional layers,  $S_a$  and  $S_{a+-}$  perform comparably in all but one layer. The two component estimator  $S_{a+-}$  is best in the dense layers.

**Image degradation** The first experiment was a direct measurement of the quality of the signal estimators of individual neurons. The second one is an indirect measurement of the quality, but it considers the whole network. We measure how the prediction (after the soft-max) for the initially selected class changes as a function of corrupting more and more patches based on the ordering assigned by the attribution (see Samek et al., 2016). This is also related to the work by Zintgraf et al. (2017). In this experiment, we split the image in non-overlapping patches of 9x9 pixels. We compute the attribution and sum all the values within a patch. We sort the patches in decreasing order based on the aggregate heat map value. In step  $n = 1..100$  we replace the first  $n$  patches with their mean per color channel to remove the information in this patch. Then, we measure how this influences the classifiers output. We use the estimators from the previous experiment to obtain the function-signal attribution heat maps for evaluation. A steeper decay indicates a better heat map.

Results are shown in Fig. 4a. The baseline, in which the patches are randomly ordered, performs worst. The linear optimized estimator  $S_a$  performs quite poorly, followed by the filter-based estimator  $S_w$ . The trivial signal estimator  $S_x$  performs just slightly better. However, the two component model  $S_{a+-}$  leads to the fastest decrease in confidence in the original prediction by a large margin. Its excellent quantitative performance is also backed up by the visualizations discussed next.

---

**Qualitative evaluation** In Fig. 4b, we compare all signal estimators on a single input image. For the trivial estimator  $S_x$ , the signal is by definition the original input image and, thus, includes the distractor. Therefore, its noisy attribution heat map shows contributions that cancel each other in the neural network. The  $S_w$  estimator captures some of the structure. The optimized estimator  $S_a$  results in slightly more structure but struggles on color information and produces dense heat maps. The two component model  $S_{a+-}$  on the right captures the original input during signal estimation and produces a crisp heat map of the attribution.

Fig. 5 shows the visualizations for six randomly selected images from ImageNet. PatternNet is able to recover a signal close to the original without having to resort to the inclusion of additional rectifiers in contrast to DeConvNet and Guided BackProp. We argue that this is due to the fact that the optimization of the pattern allows for capturing the important directions in input space. This contrasts with the commonly used methods DeConvNet, Guided BackProp, LRP and DTD, for which the correlation experiment indicates that their implicit signal estimator cannot capture the true signal in the data. Overall, the proposed approach produces the most crisp visualization in addition to being measurably better, as shown in the previous section.

**Relation to previous methods** Our method can be thought of as a generalization of the work by Haufe et al. (2014), making it applicable on deep neural networks. Remarkably, our proposed approach can solve the toy example in section 2 optimally while none of the previously published methods for deep learning are able to solve this (Bach et al., 2015; Montavon et al., 2017; Smilkov et al., 2017; Sundararajan et al., 2017; Zintgraf et al., 2017; Dabkowski & Gal, 2017; Zeiler & Fergus, 2014; Springenberg et al., 2015). Our method shares the idea that to explain a model properly one has to learn how to explain it with Zintgraf et al. (2017) and Dabkowski & Gal (2017). Furthermore, since our approach is after training just as expensive as a single back-propagation step, it can be applied in a real-time context, which is also possible for the work done by Dabkowski & Gal (2017) but not for Zintgraf et al. (2017).

## 6 CONCLUSION

Understanding and explaining nonlinear methods is an important challenge in machine learning. Algorithms for visualizing nonlinear models have emerged but theoretical contributions are scarce. We have shown that the direction of the model gradient does not necessarily provide an estimate for the signal in the data. Instead it reflects the relation between the signal direction and the distracting noise contributions (Fig. 2). This implies that popular explanation approaches for neural networks (DeConvNet, Guided BackProp, LRP) do not provide the correct explanation, even for a simple linear model. Our reasoning can be extended to nonlinear models. We have proposed an objective function for neuron-wise explanations. This can be optimized to correct the signal visualizations (PatternNet) and the decomposition methods (PatternAttribution) by taking the data distribution into account. We have demonstrated that our methods constitute a theoretical, qualitative and quantitative improvement towards understanding deep neural networks.

## ACKNOWLEDGMENTS

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement NO 657679, the BMBF for the Berlin Big Data Center BBDC (01IS14013A), a hardware donation from NVIDIA. We thank Sander Dieleman, Jonas Degraeve, Ira Korshunova, Stefan Chmiela, Malte Esders, Sarah Hooker, Vincent Vanhoucke for their comments to improve this manuscript. We are grateful to Chris Olah and Grégoire Montavon for the valuable discussions.

## REFERENCES

- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.

- 
- David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(Jun):1803–1831, 2010.
- James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: A cpu and gpu math compiler in python. In *Proc. 9th Python in Science Conf*, pp. 1–7, 2010.
- Piotr Dabkowski and Yarin Gal. Real time image saliency for black box classifiers. In *NIPS 2017*, 2017.
- Sander Dieleman, Jan Schlüter, Colin Raffel, Eben Olson, Søren Kaae Sønderby, Daniel Nouri, Daniel Maturana, Martin Thoma, Eric Battenberg, Jack Kelly, et al. Lasagne: First release. *Zenodo: Geneva, Switzerland*, 2015.
- Stefan Haufe, Frank Meinecke, Kai Görzen, Sven Dähne, John-Dylan Haynes, Benjamin Blankertz, and Felix Bießmann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage*, 87:96–110, 2014.
- Pieter-Jan Kindermans, Kristof Schütt, Klaus-Robert Müller, and Sven Dähne. Investigating the influence of noise and distractors on the interpretation of neural networks. *arXiv preprint arXiv:1611.07270*, 2016.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. doi: 10.1038/nature14539.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
- Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems*, pp. 3387–3395, 2016.
- 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, 2015.
- Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE Transactions on Neural Networks and Learning Systems*, 2016. doi: 10.1109/TNNLS.2016.2599820.
- Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015.
- Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *CoRR*, abs/1605.01713, 2016. URL <http://arxiv.org/abs/1605.01713>.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi: 10.1038/nature16961.

- 
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR*, 2014.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. In *ICLR*, 2015.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *ICML 2017*, 2017.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pp. 3104–3112, 2014.
- Jason Yosinski, Jeff Clune, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. In *ICML Workshop on Deep Learning*, 2015.
- Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pp. 818–833. Springer, 2014.
- Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. Visualizing deep neural network decisions: Prediction difference analysis. In *ICLR*, 2017.