

# Reasoning and Provenance on Neural Networks

Antoine Groudiev

L3, DI, ENS

Silviu Maniu

Slide Team, LIG, UGA

14th June 2024

## Abstract

Recently, neural networks allowed computers to solve numerous problems from diverse machine learning fields, such as natural language processing and computer vision. Compared to traditional algorithms, machine learning models have proven both more successful and more difficult to interpret. Neural networks are considered as black boxes unable to easily explain themselves, that is justifying the reasons that led them to make a prediction. Layer-wise Relevance Propagation (LRP) is a technique that has been introduced to provide explainability by identifying the input features relevant to the output choice. In parallel, research in the databases field developed annotations techniques to compute provenance for queries. In this paper, we extend LRP propagation rules to semiring-based provenance annotations of the network, for LRP to gain in expressivity.

## 1 Introduction

### 1.1 Problem statement

Deep neural networks have proven successful for solving with high accuracy machine learning problems. The expressivity of the class of functions generated by neural networks, combined with the relative simplicity of their training, make such models versatile tools to learn the relationship between the inputs and outputs of a dataset.

However, this versatility comes at the cost of poor interpretability: a neural network simply represents a function from one high-dimensional space to another, but provides no justification nor explanation for a given execution. If metrics such as the accuracy over a testing set provide confidence in the fact that the model is able to correctly classify inputs similar to the training set, no guarantee is given that the model generalizes well. Real-world examples show that networks can overfit the input data, or even take shortcuts instead of learning the intended solution [2]. For the user to have confidence in its predictions, a neural network should therefore be able to highlight the patterns in the input data that it actually learned.

### 1.2 Layer-wise Relevance Propagation

Layer-wise Relevance Propagation (LRP) [1] has been introduced as a technique to explain an execution of a neural network. LRP is a procedure propagating the output of the function backward in the network, using diverse rules to compute the *relevance* of a neuron depending on the relevances of the neurons of the upstream layer. LRP introduces the notion of *relevance score* for a neuron, intuitively quantifying the contribution of this neuron to the classification of final classification. A high relevance score indicates that the neuron led to the activation of the considered output; a negative relevance score represent neurons that increased the activation of another output neuron instead of the one considered.

### 1.2.1 Setup and notations

In the following, we consider a deep neural network used for a classification task. We assume that it uses the rectifier activation function<sup>1</sup>, which is the case in most applications. To ease the notation, we will not consider biases but instead assume that the first neuron of each layer represents the bias.

Let  $L$  be the number of layers of the network. We denote by  $(a_k^{(l)})_k$  the activations of the network. Notably,  $(a_k^{(1)})_k$  is the input data, and  $(a_k^{(L)})_k$  is the output prediction. We denote by  $(w_{j,k}^{(l)})_{j,k}$  the weights connecting the  $l$ -th layer to the  $(l+1)$ -th layer. To simulate the weights, we set:

$$\forall l \in \llbracket 1, L \rrbracket, \quad a_0^{(l)} = 1$$

and we define  $w_{0,k}^{(l)}$  to be the bias of the  $k$ -th neuron of the  $(l+1)$ -th layer. The forward propagation rule of a deep rectifier network is therefore:

$$\forall l \in \llbracket 1, L-1 \rrbracket, \forall k, \quad a_k^{(l+1)} = \text{ReLU} \left( \sum_{j=0} a_j^{(l)} w_{j,k} \right) = \max \left( 0, \sum_{j=0} a_j^{(l)} w_{j,k} \right) \quad (1.2.1)$$

We denote by  $R_j^{(l)}$  the relevance of the  $j$ -th neuron of the  $l$ -th layer. We assume that the output layer represents a one-hot encoding, that is that the belonging of the input to the  $i$ -th class is represented by an output vector is of the form  $(0, \dots, a_i^{(L)}, \dots, 0)$ , that where the only non-null coefficient is in the  $i$ -th position. Finally, we denote by  $y$  the label of a classified input. To the label  $y = i$  is associated the output vector  $(0, \dots, a_i^{(L)}, \dots, 0)$ .

### 1.2.2 Propagation rules

Relevance scores are initialized for the output layer, and are set to the output activation for the correct class, that is:

$$R_i^{(L)} = \begin{cases} a_i^{(L)} & \text{if } i = y \\ 0 & \text{otherwise} \end{cases} \quad (1.2.2)$$

The simplest LRP rule is called LRP-0. It propagates the relevance to a neuron of the lower layer proportionally to its contribution to each of the neuron of the next layer:

$$R_j^{(l)} = \sum_k \frac{a_j^{(l)} w_{j,k}}{\sum_{j'} a_{j'}^{(l)} w_{j',k}} R_k^{(l+1)} \quad (1.2.3)$$

The denominator  $\sum_{j'} a_{j'}^{(l)} w_{j',k}$  guarantees a conservation property, that is that for any layer  $l$ :

$$\sum_j R_j^{(l)} = \sum_k R_k^{(l+1)}$$

This allows to keep the information regarding the total activation of the final layer.

The application of this simple rule can lead into noisy results that do not scale well. An overview of variations of LRP-0 is provided by [3]; not all rules are suitable for all layers. Complex deep neural networks architectures benefit from enhanced rules such as LRP- $\varepsilon$  or LRP- $\gamma$ , which provide more stable explanations.

---

<sup>1</sup>That is  $\text{ReLU}(x) = \max(0, x)$ , where ReLU stands for *Rectified Linear Unit*.

Notably, the input layer must be handled using a different rule, since it does not receive its input from ReLU activations, but directly from the input data. The  $z^{\mathcal{B}}$  rule is used in [3] to propagate from layer 2 to 1 (input layer):

$$R_j^{(1)} = \sum_k \frac{x_k w_{j,k} - l_j w_{j,k}^+ - h_j w_{j,k}^-}{\sum_{j'} x_k w_{j',k} - l_j w_{j',k}^+ - h_j w_{j',k}^-} R_k^{(2)} \quad (1.2.4)$$

where  $(\cdot)^+ = \max(0, \cdot)$ ,  $(\cdot)^- = \min(0, \cdot)$ . The parameters  $l_i$  and  $h_i$  respectively define the theoretical minimum and maximum values of the inputs  $x_i$ ; for instance we might have  $l_i = 0$  and  $h_i = 255$  for pixels over 8 bits.

### 1.3 Semiring-based provenance annotations

## References

- [1] Sebastian Bach et al. ‘On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation’. In: *PLOS ONE* (2015), pp. 1–46. DOI: [10.1371/journal.pone.0130140](https://doi.org/10.1371/journal.pone.0130140). URL: <https://doi.org/10.1371/journal.pone.0130140>.
- [2] Robert Geirhos et al. ‘Shortcut learning in deep neural networks’. In: *Nature Machine Intelligence* 2 (2020), pp. 665–673.
- [3] Grégoire Montavon et al. ‘Layer-Wise Relevance Propagation: An Overview’. In: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer International Publishing, 2019, pp. 193–209. URL: [https://doi.org/10.1007/978-3-030-28954-6\\_10](https://doi.org/10.1007/978-3-030-28954-6_10).
- [4] Yann Ramusat, Silviu Maniu and Pierre Senellart. ‘Provenance-Based Algorithms for Rich Queries over Graph Databases’. In: *EDBT 2021 - 24th International Conference on Extending Database Technology*. 2021. URL: <https://inria.hal.science/hal-03140067>.