

Extending Layerwise Relevance Propagation using Semiring Annotations

Antoine Groudiev
L3, ENS Ulm

Silviu Maniu – Supervisor
SLIDE Team, LIG

September 2, 2024



Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Semiring generalization of the LRP rule

Results over the MNIST dataset

Applications

Image mask computation

Network pruning using LRP ranking

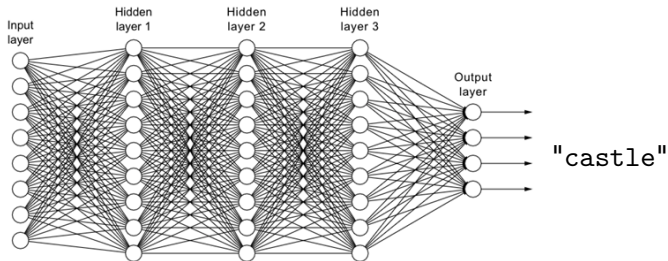
Comparison to image perturbation

Handling CNNs

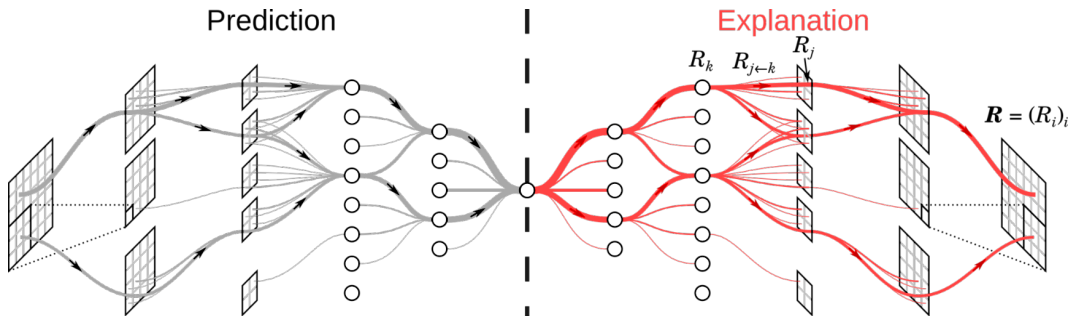
Computing relevance for convolutional layers

Results for the VGG-16 network

Problem statement



Layerwise Relevance Propagation [11]



Layerwise Relevance Propagation

Initialization

Initialization:

$$R_i^{(L)} = \begin{cases} a_i^{(L)} & \text{if } i = y \text{ (the class we want)} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

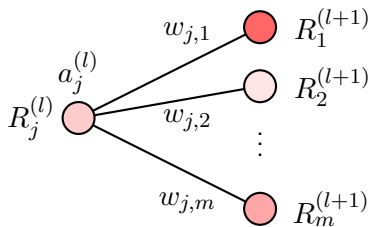
$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 4.2 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \rightarrow \text{"goldfish"} \\ \rightarrow \text{"street sign"} \\ \\ \rightarrow \text{"castle"} \\ \\ \rightarrow \text{"printer"} \end{matrix}$$

Layerwise Relevance Propagation

Propagation

LRP-0 rule:

$$R_j^{(l)} = \sum_k \frac{a_j^{(l)} w_{j,k}}{\sum_{j'} a_{j'}^{(l)} w_{j',k}} \cdot R_k^{(l+1)} \quad (\text{LRP-0})$$



Other rules exist (LRP- ϵ , LRP- γ , z^B)

LRP Results visualization

VGG-16 on ImageNet dataset



Figure 2: Reference image

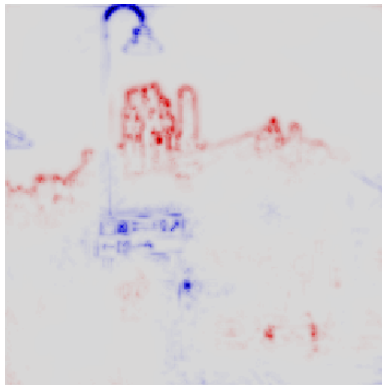


Figure 3: Relevance for the class "castle"

Semiring-based provenance annotations [7, 12]

Definition (Semiring)

A semiring $(\mathbb{K}, \oplus, \otimes, 0, 1)$ is such that:

- \otimes distributes over \oplus ,
- $(\mathbb{K}, \oplus, 0)$ is a commutative monoid,
- $(\mathbb{K}, \otimes, 1)$ is a monoid such that 0 is absorbing

Example

The following structures are semirings:

- Real semiring: $(\mathbb{R}, +, \times, 0, 1)$
- Boolean semiring: $(\{\perp, \top\}, \vee, \wedge, \perp, \top)$
- Counting semiring: $(\mathbb{N}, +, \times, 0, 1)$
- Viterbi semiring: $([0, 1], \max, \times, 0, 1)$

Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Semiring generalization of the LRP rule

Results over the MNIST dataset

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Handling CNNs

Computing relevance for convolutional layers

Results for the VGG-16 network

Semiring generalization of the LRP rule

Consider a semiring $(\mathbb{K}, \oplus, \otimes, \mathbb{0}, \mathbb{1})$

Conversion function:

$$\Theta : \mathbb{R} \longrightarrow \mathbb{K}$$

Initialization:

$$R_i^{(L)} = \begin{cases} \mathbb{1} & \text{if } i = y \\ \mathbb{0} & \text{otherwise} \end{cases} \quad (2)$$

Propagation rule:

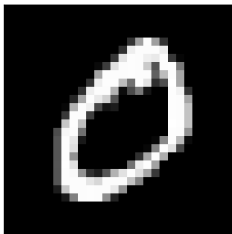
$$R_j^{(l)} = \bigoplus_k \Theta \left(\frac{a_j^{(l)} w_{j,k}}{\sum_{j'} a_{j'}^{(l)} w_{j',k}} \right) \otimes R_k^{(l+1)} \quad (\mathbb{K}\text{-LRP})$$

Boolean Semiring

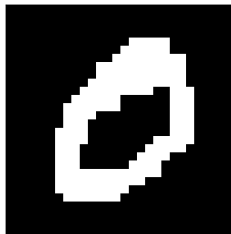
$(\{\perp, \top\}, \vee, \wedge, \perp, \top)$

$$\Theta = x \mapsto \begin{cases} \top & \text{if } x \geq \theta \\ \perp & \text{otherwise} \end{cases}$$

Reference



Boolean Semiring

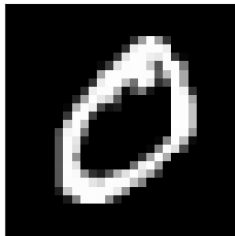


Counting Semiring

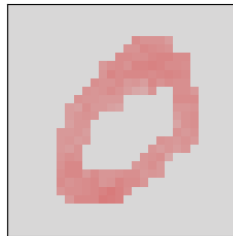
$(\mathbb{N}, +, \times, 0, 1)$

$$\Theta = x \mapsto \begin{cases} 1 & \text{if } x \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

Reference



Counting semiring

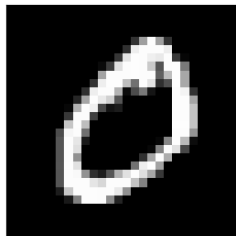


Viterbi Semiring

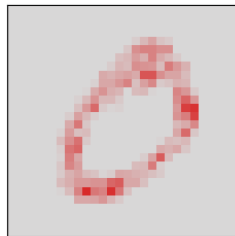
$([0, 1], \max, \times, 0, 1)$

$$R_j^{(l)} = \max_k \underbrace{\left(\frac{|a_j^{(l)} w_{j,k}^{(l)}|}{\max_{j'} |a_{j'}^{(l)} w_{j',k}^{(l)}|} \right)}_{\in [0,1]} \cdot R_k^{(l+1)}$$

Reference



Viterbi semiring



Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Semiring generalization of the LRP rule

Results over the MNIST dataset

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

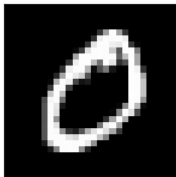
Handling CNNs

Computing relevance for convolutional layers

Results for the VGG-16 network

Class-wise mask – Boolean semiring

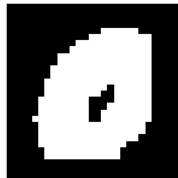
Reference



Class-wise AND (5 examples)



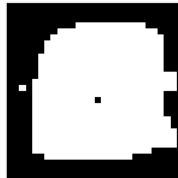
Class-wise OR
(5 examples)



Class-wise AND
(100 examples)

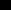
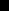


Class-wise OR
(100 examples)

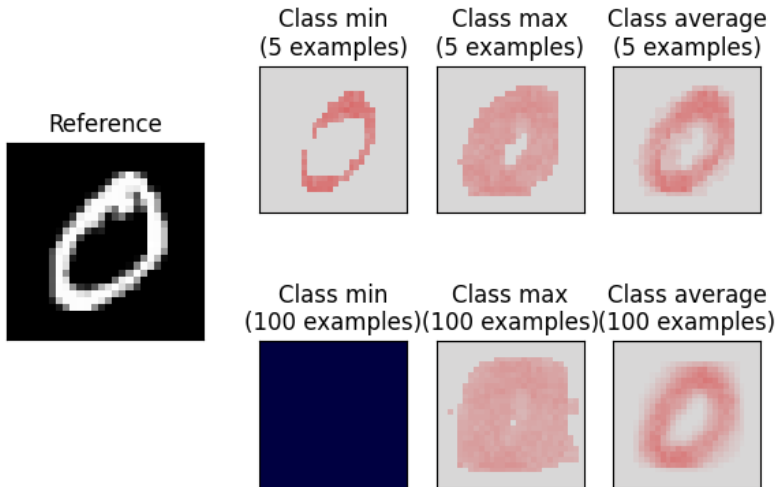


All classes mask – Boolean semiring

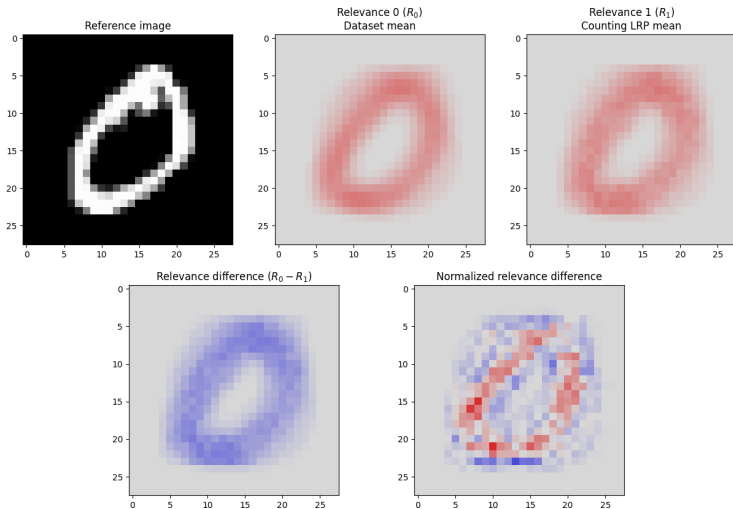
1. *Journal of Management Studies*, 1996, 33, 1, 1-14.

[illegible]

Class-wise mask – Counting semiring



Comparison to dataset mean



Network pruning using LRP ranking

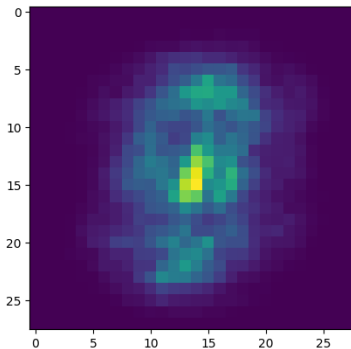


Figure 4: Relevance mean over the training dataset
(Input layer)

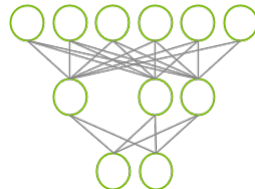
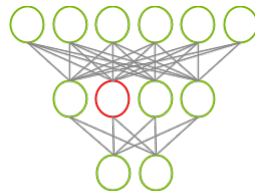


Figure 1 is a line graph showing the relationship between the percentage of pruned parameters (x-axis, 0% to 100%) and the loss (y-axis, 0.0 to 2.0). The graph compares five different pruning methods: a baseline (black line), a method with a higher loss (orange line), a method with a moderate loss (blue line), a method with a lower loss (red line), and the proposed method (green line). The proposed method (green line) consistently shows the lowest loss across the range of pruned parameters, particularly for higher percentages of pruning (above 60%).

Pruned parameters (%)	Baseline (Black)	Orange Line	Blue Line	Red Line	Proposed (Green)
0	0.15	0.15	0.15	0.15	0.15
20	0.15	0.15	0.15	0.15	0.15
40	0.25	0.30	0.20	0.20	0.25
60	0.60	0.80	0.40	0.30	0.60
80	1.30	1.50	0.80	0.60	1.30
100	2.30	2.30	2.00	2.00	2.30

Figure 1 is a line graph showing the relationship between the percentage of pruned parameters (x-axis, 0% to 100%) and the resulting accuracy (y-axis, 20% to 80%) for six different LRP methods. The methods are: Vieri LRP (blue line), Boolean LRP (orange line), counting LRP (green line), Classic LRP (red line), ℓ_2 pruning (purple line), and ℓ_1 pruning (brown line). All methods start with an accuracy of approximately 95% at 0% pruned parameters. As the percentage of pruned parameters increases, the accuracy for all methods decreases. Vieri LRP consistently maintains the highest accuracy across the range of pruned parameters, while ℓ_1 pruning maintains the lowest accuracy. The accuracy for all methods drops sharply as the percentage of pruned parameters approaches 100%.

Pruned parameters (%)	Vieri LRP (%)	Boolean LRP (%)	counting LRP (%)	Classic LRP (%)	ℓ_2 pruning (%)	ℓ_1 pruning (%)
0	95	95	95	95	95	95
20	95	95	95	95	95	95
40	95	95	95	95	95	95
60	95	90	85	95	95	95
80	90	55	50	85	80	80
100	15	10	10	10	10	10

Comparison to image perturbation [5]

Accuracies per attack zone
Kernel size: 4 — Step: 1

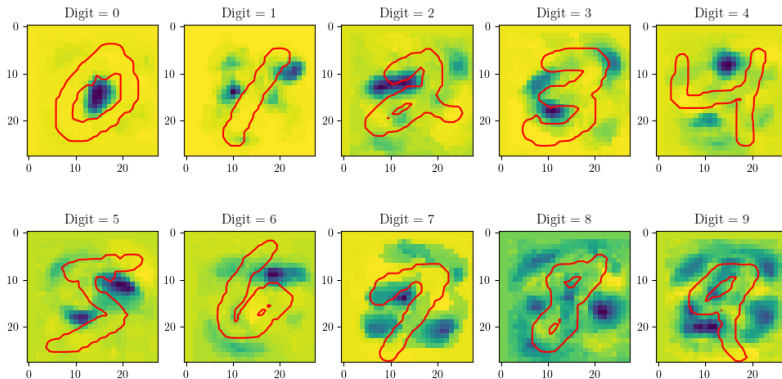


Figure 5: Accuracies per attack zone

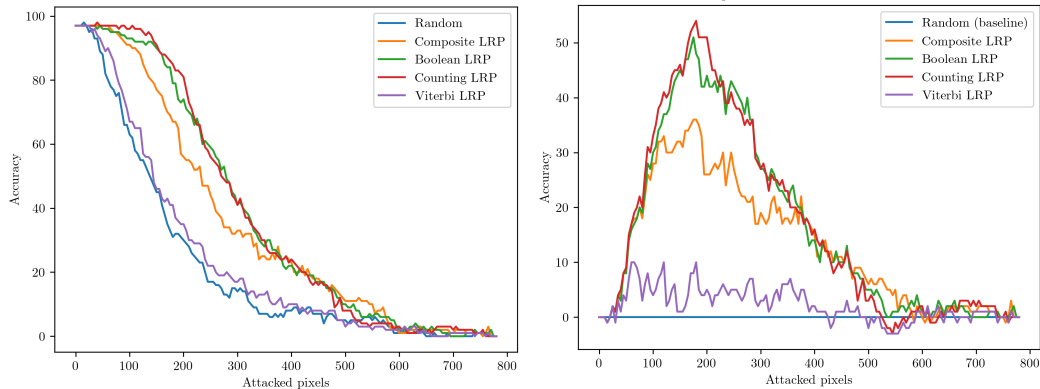


Figure 6: Accuracy drop for multiple pixels attacks strategies.

Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Semiring generalization of the LRP rule

Results over the MNIST dataset

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Handling CNNs

Computing relevance for convolutional layers

Results for the VGG-16 network

Computing relevance for convolutional layers

$$R_j^{(l)} = \underbrace{\bigoplus_k \underbrace{\Theta \left(\frac{a_j^{(l)} w_{j,k}}{\sum_{j'} a_{j'}^{(l)} w_{j',k}} \right)}_{\text{Convolution over layer } l}}_{\text{Convolution over layer } l+1} \otimes R_k^{(l+1)} \quad (\mathbb{K}\text{-LRP})$$

VGG-16 network

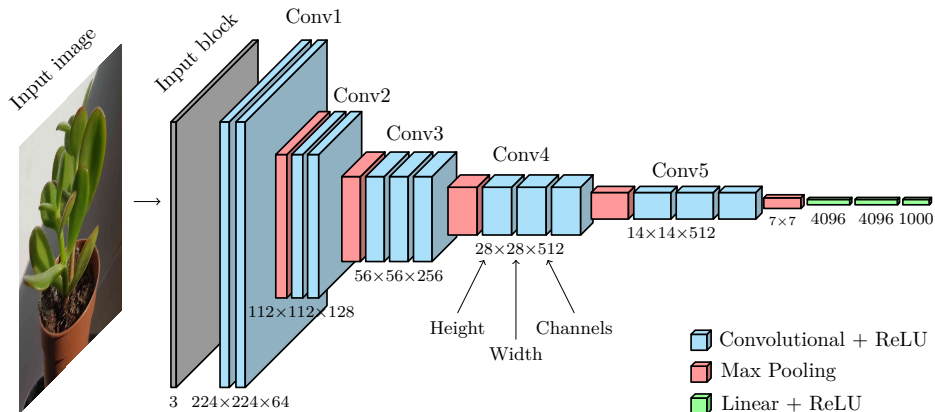
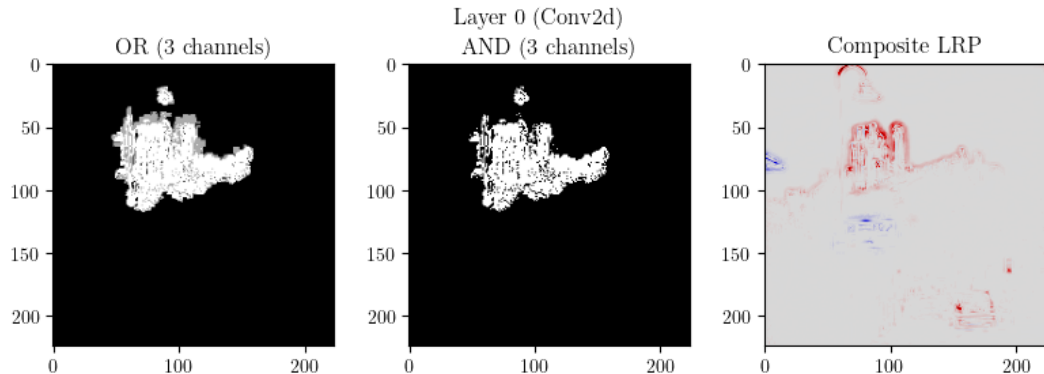


Figure 7: Architecture of the VGG-16 network.

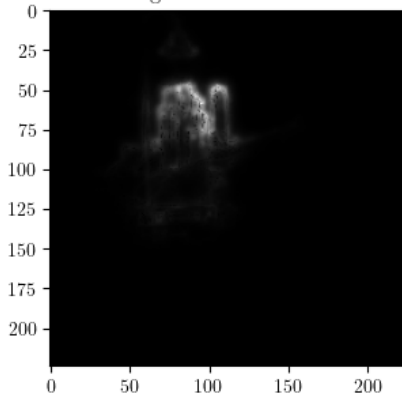
Results for VGG-16: Boolean semiring



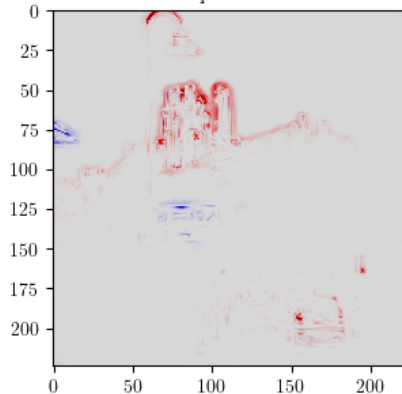
Results for VGG-16: Counting semiring

Layer 0 (Conv2d)

Counting - Sum over 3 channels



Composite LRP



- [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. URL: <https://doi.org/10.1371/journal.pone.0130140>.
- [2] Nick Cammarata et al. "Thread: Circuits". In: *Distill* (2020). <https://distill.pub/2020/circuits>. DOI: 10.23915/distill.00024.
- [3] Hongrong Cheng, Miao Zhang, and Javen Qinfeng Shi. "A Survey on Deep Neural Network Pruning-Taxonomy, Comparison, Analysis, and Recommendations". In: (2023). URL: <https://arxiv.org/abs/2308.06767>.
- [4] Marina Danilevsky et al. "A survey of the state of explainable AI for natural language processing". In: *arXiv preprint* (2020). URL: <https://arxiv.org/abs/2010.00711>.
- [5] Ruth C Fong and Andrea Vedaldi. "Interpretable explanations of black boxes by meaningful perturbation". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 3429–3437. URL: <https://arxiv.org/abs/1704.03296>.
- [6] Robert Geirhos et al. "Shortcut learning in deep neural networks". In: *Nature Machine Intelligence* 2 (2020), pp. 665–673. URL: <https://arxiv.org/abs/2004.07780>.

- [13] Wojciech Samek et al. “Evaluating the visualization of what a deep neural network has learned”. In: *IEEE transactions on neural networks and learning systems* 28.11 (2016), pp. 2660–2673. URL: <https://arxiv.org/abs/1509.06321>.
- [14] Ramprasaath R Selvaraju et al. “Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 618–626. URL: <https://arxiv.org/pdf/1610.02391>.
- [15] Pierre Senellart. “Provenance and probabilities in relational databases”. In: *ACM SIGMOD Record* 46.4 (2018), pp. 5–15. URL: <https://inria.hal.science/hal-01672566>.
- [16] Pierre Senellart et al. “ProvSQL: Provenance and probability management in postgresql”. In: *Proceedings of the VLDB Endowment (PVLDB)* (2018). URL: <https://inria.hal.science/hal-01851538>.
- [17] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. 2015. URL: <https://arxiv.org/abs/1409.1556>.
- [18] Matthew D Zeiler and Rob Fergus. “Visualizing and understanding convolutional networks”. In: *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I* 13. Springer. 2014, pp. 818–833. URL: <https://arxiv.org/abs/1311.2901>.